

Reward Optimizing Recommendation with Deep Learning and Maximum Inner Product Search

ECML PKDD Tutorial 2022

Imad Aouali, Amine Benhalloum, Martin Bompaire,
Achraf Ait Sidi Hammou, Sergey Ivanov, Benjamin
Heymann, David Rohde, Otmane Sakhi, Flavian
Vasile, Maxime Vono

Agenda

1

Introduction
Practice → Science

2

Framework
Contextual Bandits

3

**Single Item
Recommendation**

4

**Slate
Recommendation**

5

Conclusion

Introduction

Practice

Introduction (Practice)

What are we trying to do?

Recommend the **best** items to a **user**

- **Recommend:** Within the context of online advertising, low latency constraints, slates
- **Best:** Maximizing a given reward (clicks, sales, ...)
- **Items:** Tens of millions of products
- **User:** History of past user interactions, adapt to change in behaviour quickly, users enter and leave the system continually

Introduction (Practice)

General Framework

We are interested in decision rules of the form: $\text{score}(\mathbf{u}, a) \sim \mathbf{u}^T \boldsymbol{\beta}_a$

\mathbf{u} - The user embedding of dimension D

$\boldsymbol{\beta}_a$ - The embedding of item a of dimension D

$\boldsymbol{\beta}$ - The matrix of item embeddings of dimension $D \times P$

Introduction (Practice)

General Framework

We are interested in decision rules of the form: $\text{score}(\mathbf{u}, a) \sim \mathbf{u}^T \boldsymbol{\beta}_a$

Why? To find the best action we need:

$$a^* = \operatorname{argmax}_{a' \in \{1, \dots, P\}} \text{score}(\mathbf{u}, a')$$

If P is large, then this is very slow for general score functions.

Using Dot Product enables us to leverage (Approximate) Maximum Inner Product Search techniques.

Introduction (Practice)

General Framework

We are interested in decision rules of the form: $\text{score}(\mathbf{x}, a) \sim f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_a$

\mathbf{x} - The **context** which is the information we use to personalize on.

f_{Ξ} - A mapping from context to user embedding.

Introduction (Practice)

General Framework

We are interested in decision rules of the form: $\text{score}(\mathbf{x}, a) \sim \mathbf{f}_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_a$

Why?

- Users enter and leave the system constantly.
- We want to adapt to change in behavior.
- Generalize across users.

Introduction (Practice)

General Framework

We are interested in decision rules of the form: $\text{score}(\mathbf{x}, \mathbf{a}) \sim f_{\Xi}(\mathbf{x})^T h_{\beta}(\mathbf{a})$

\mathbf{a} - Set of initial pre-trained embeddings, product attributes and features ...

Introduction (Practice)

General Framework

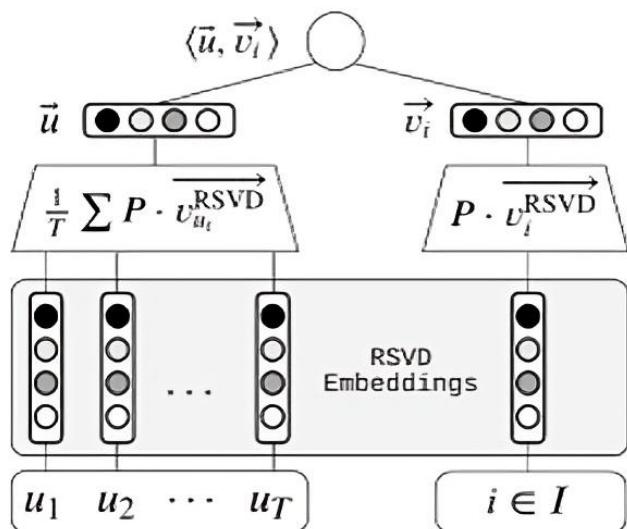
We are interested in decision rules of the form: $\text{score}(\mathbf{x}, \mathbf{a}) \sim f_{\Xi}(\mathbf{x})^T h_{\beta}(\mathbf{a})$

Why ?

- Learn from vast amounts of heterogeneous sources of data, side information, etc.
- Enable transfer between tasks.

Introduction (Practice)

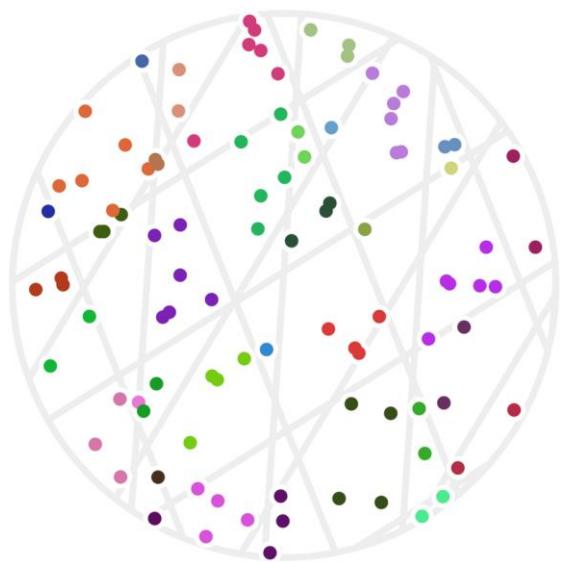
General Framework



- The now ubiquitous Two Tower style model ([Criteo DeepKNN](#), [Facebook DLRM](#) ...)
- You can make the towers as shallow or deep as you want/need #deeplearning.

Introduction (Practice)

General Framework



$$\text{argsort}(f_{\Xi}(x)^T \beta)_{1:K}$$

where

$$f_{\Xi}(x) = D$$

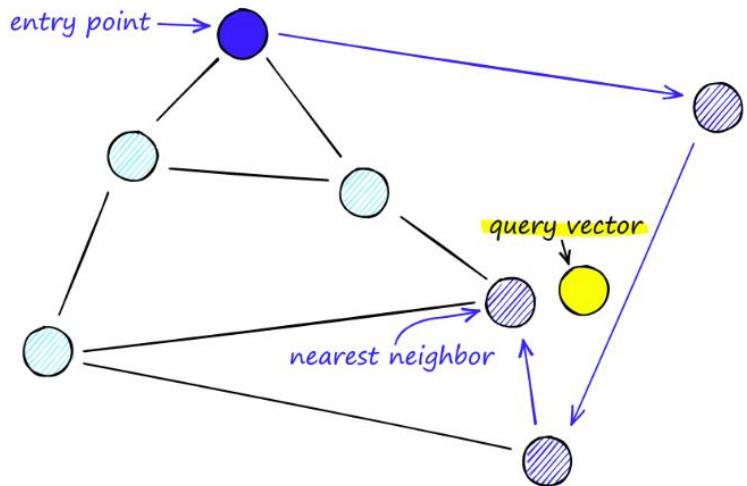
$$\beta = D \times P$$

- Scores computation is $O(P \cdot D)$
- Exact sorting is $O(P \log P)$
- Exact top- K is $O(P)$

Solution: Approximate algorithms that “look like” $O(\log P)$
e.g., Navigable Small World

Introduction (Practice)

Side note: Hierarchical Navigable Small Worlds (HSNW)



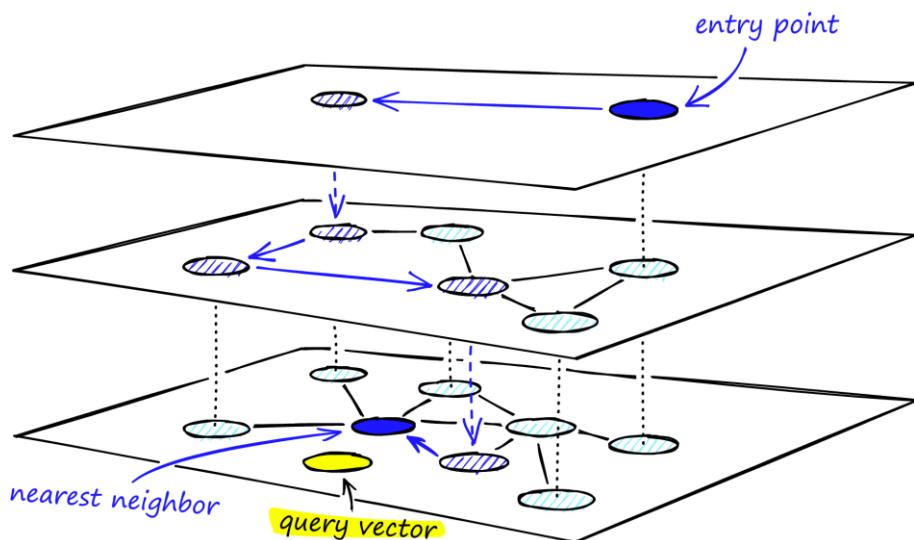
Credits: <https://www.pinecone.io/learn/hnsw/>

General recipe for Graph Based NN methods:

- Build a graph where close vectors are linked together.
- At query time, route your query vector through the graph greedily, navigating to the closest neighbor at each time.
- Stop when you hit a local minimum.

Introduction (Practice)

Side note: Hierarchical Navigable Small Worlds (HSNW)



Credits: <https://www.pinecone.io/learn/hnsw/>

The hierarchical part

- We start by a **zoomed out** view of the neighborhood graph
- And we zoom in after routing at each successive layer
- More fine grained search as we go deeper

Introduction (Practice)

Python Tutorial 1 (Exact vs. Approximate MIPS)

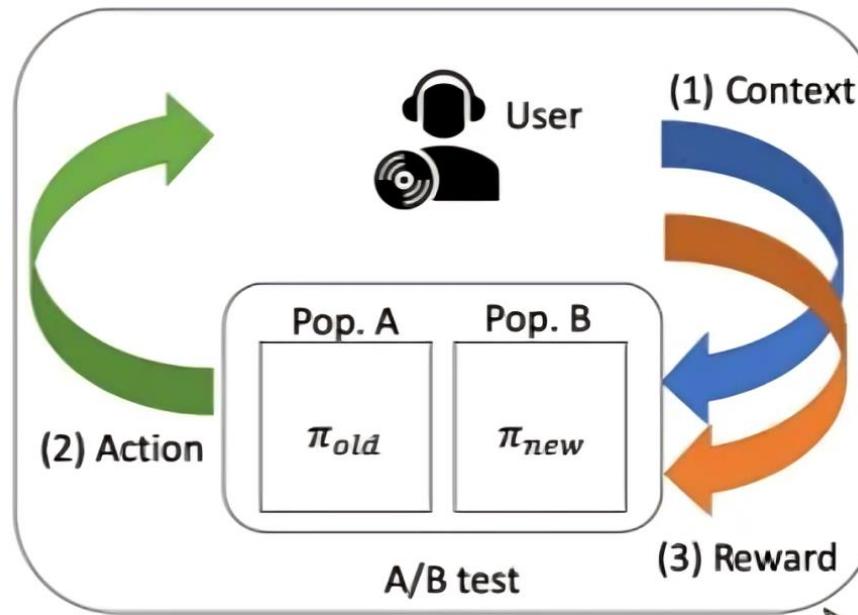
[https://colab.research.google.com/github/otmhi/Reward-Optimizing-Reco/blob/main/ECML_Rewards\[...\]izing_Slate_Recommendation_with_DL_and_MIPS_Part_1.ipynb](https://colab.research.google.com/github/otmhi/Reward-Optimizing-Reco/blob/main/ECML_Rewards[...]izing_Slate_Recommendation_with_DL_and_MIPS_Part_1.ipynb)

Introduction

Science

Introduction (Science)

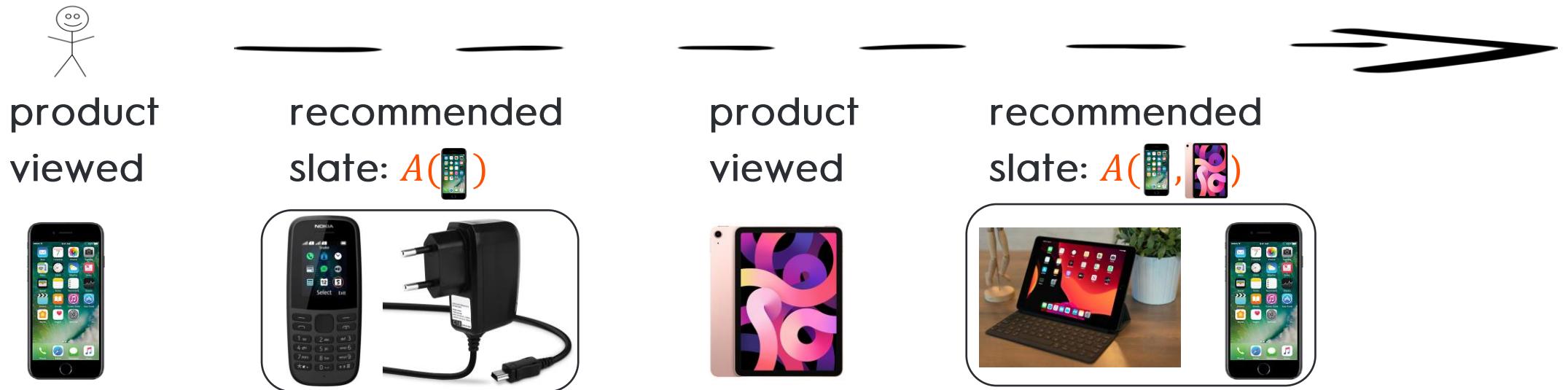
Reward Optimizing Recommendation



Introduction (Science)

Reward Optimizing Recommendation

Population A



$A(\cdot)$ is the algorithm we test on population A: it is a mapping from a list of viewed items to a list (slate) of recommended items.

Introduction (Science)

Reward Optimizing Recommendation

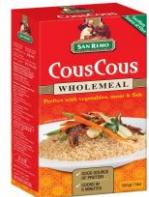
Population B



product
viewed

product
viewed

recommended
slate: $B(\cdot)$



product
viewed

recommended
slate: $S(\cdot)$



$B(\cdot)$ is the algorithm we test on population B: it is a mapping from a list of viewed items to a list (or slate) of recommended items.

Introduction (Science)

Reward Optimizing Recommendation

Which is better $A(\cdot)$ or $B(\cdot)$?

- Recall that $A(\cdot)$ and $B(\cdot)$ are mappings from lists of items to lists of items e.g.

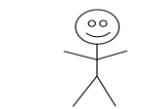
$$\text{[CousCous Box, Gold Beer, iPhone]} = B(\text{[CousCous Box, DAAWAT Brown Rice]})$$

- In general, $A(\cdot)$ and $B(\cdot)$ are constrained by engineering limitations.
- An A/B test can measure the timeline reward accurately (If we make the stable unit treatment value assumption).

Introduction (Science)

Reward Optimizing Recommendation

How do we do this offline?



product
viewed



product
viewed



recommended
slate: $B(\cdot)$



product
viewed



recommended
slate: $B(\cdot)$



Introduction (Science)

Reward Optimizing Recommendation

How do we do this offline?



product
viewed



product
viewed



recommended
slate: $B(\cdot)$



product
viewed



recommended
slate: $B(\cdot)$



→ It is really difficult !

Introduction (Science)

Non-Reward Optimizing Recommendation with Pseudo Rewards

Our problem is that our recommender is a mapping:

$$\text{Food Box} \times \text{Beer Bottle} \times \text{Smartphone} = B(\text{Food Box}, \text{Food Box})$$

And we're trying to answer the counterfactual question: **what would happen if I had a different mapping?** from data generated by a previous mapping where rewards could be delayed...

Introduction (Science)

Non-Reward Optimizing Recommendation with Pseudo Rewards

How we do this offline?



product
viewed



Introduction (Science)

Non-Reward Optimizing Recommendation with Pseudo Rewards

How we do this offline?



product
viewed



Introduction (Science)

Non-Reward Optimizing Recommendation with Pseudo Rewards

How we do this offline?



product
viewed

product
viewed

product
viewed



Becomes a positive example

Introduction (Science)

Non-Reward Optimizing Recommendation with Pseudo Rewards

How we do this offline?



product
viewed

product
viewed

product
viewed



$$L(\Xi, \beta) = \log \sigma\{f_{\Xi}(v_1 = \text{[red box]}, v_2 = \text{[red box}})^T \beta_a^+ - f_{\Xi}(v_1 = \text{[green box]}, v_2 = \text{[green box}})^T \beta_a^-\}$$

Becomes a positive example

Introduction (Science)

Non-Reward Optimizing Recommendation with Pseudo Rewards

How we do this offline?



product
viewed

product
viewed

product
viewed



$$L(\Xi, \beta) = \log \sigma\{f_{\Xi}(v_1 = \text{[red box]}, v_2 = \text{[red box}})^T \beta_a^+ - f_{\Xi}(v_1 = \text{[green box]}, v_2 = \text{[green box}})^T \beta_a^-\}$$

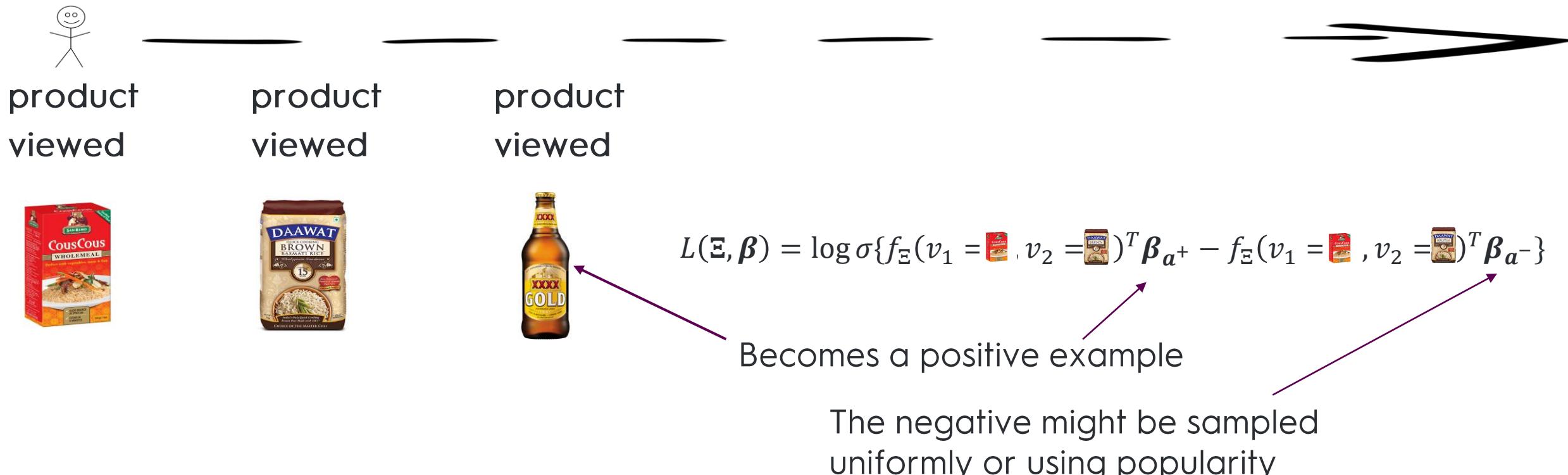
Becomes a positive example

The negative might be sampled
uniformly

Introduction (Science)

Non-Reward Optimizing Recommendation with Pseudo Rewards

How we do this offline?



Introduction (Science)

Non-Reward Optimizing Recommendation with Pseudo Rewards

How do we use recommendation feedback ?

First answer, use a ranking loss function while training on historical recommendations and feedback.

e.g., Some pairwise ranking loss such as BPR

Introduction (Science)

Non-Reward Optimizing Recommendation with Pseudo Rewards

How do we use recommendation feedback ?



product
viewed



Introduction (Science)

Non-Reward Optimizing Recommendation with Pseudo Rewards

How do we use recommendation feedback ?



product
viewed



Introduction (Science)

Reward Optimizing Recommendation

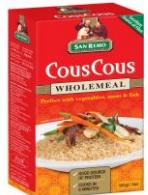
How do we use recommendation feedback ?



product
viewed

product
viewed

recommended
slate:



No click

click

No click

Introduction (Science)

Reward Optimizing Recommendation

How do we use recommendation feedback ?



product
viewed

product
viewed

recommended
slate:



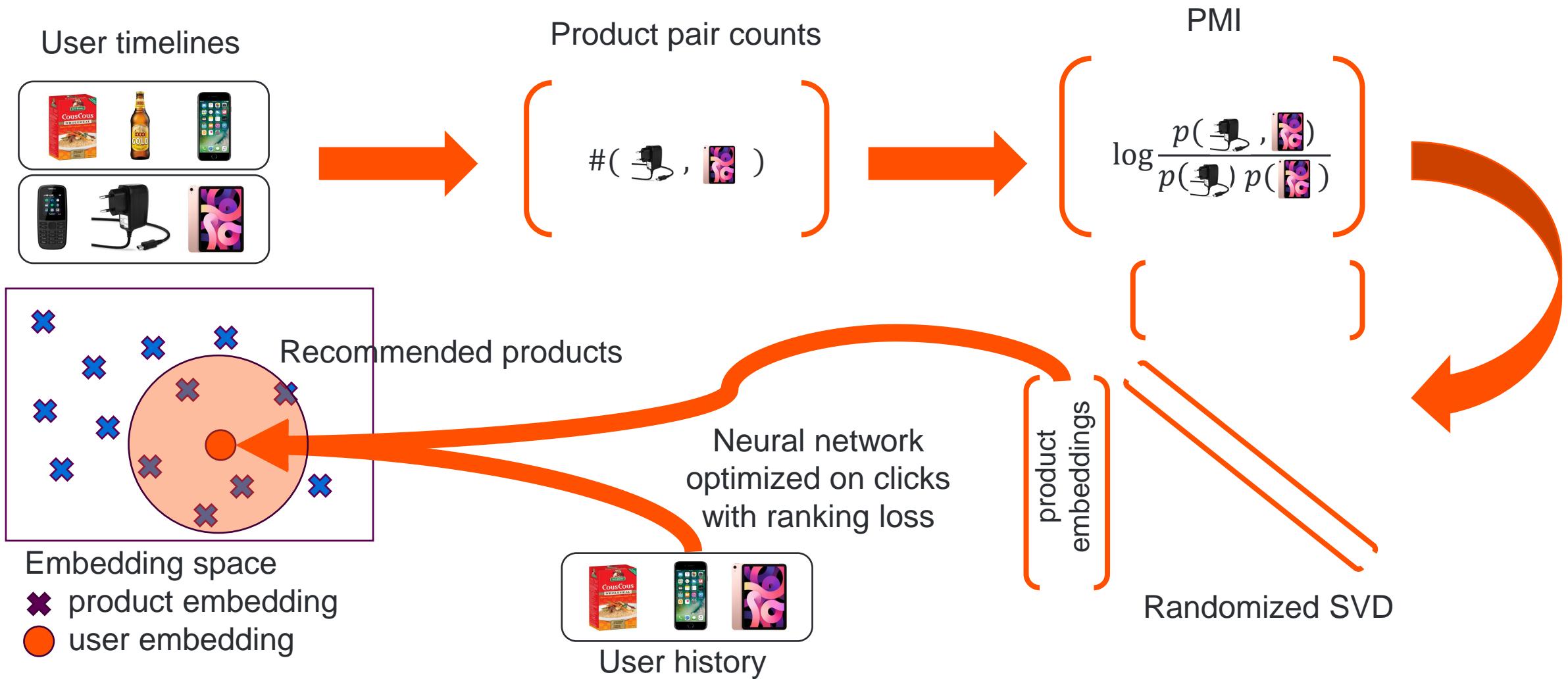
No click

click

No click

$$L(\Xi, \beta) = \log \sigma\{f_{\Xi}(v_1 = \text{rice}, v_2 = \text{beer})^T \beta_{a^+} - f_{\Xi}(v_1 = \text{rice}, v_2 = \text{snack})^T \beta_{a^-}\}$$

Recommendation at Criteo

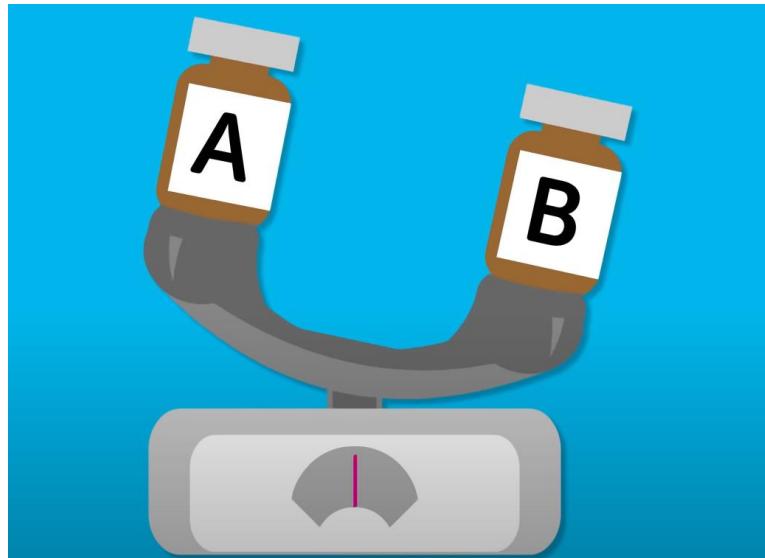


Introduction

Reward on the Timeline

Introduction (Reward on the Timeline)

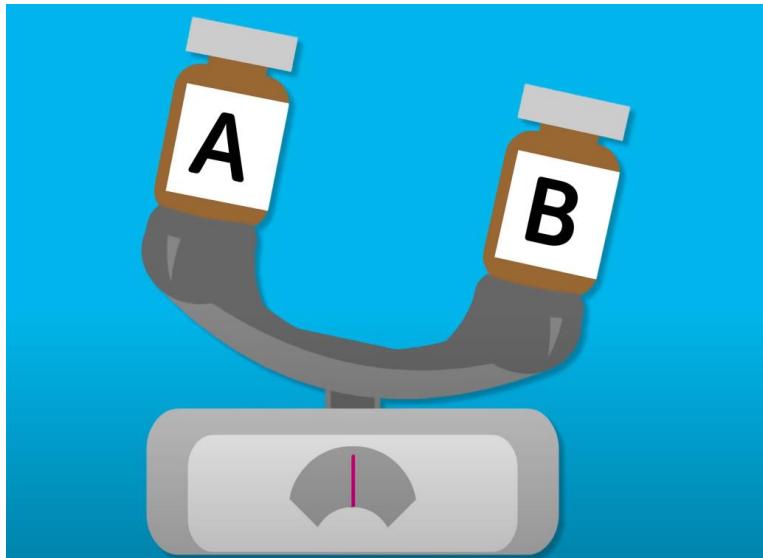
Randomized Control Trials (RCT) in Medicine



In medicine a randomized control trial, splits the population into two groups.

Introduction (Reward on the Timeline)

Randomized Control Trials (RCT) in Medicine

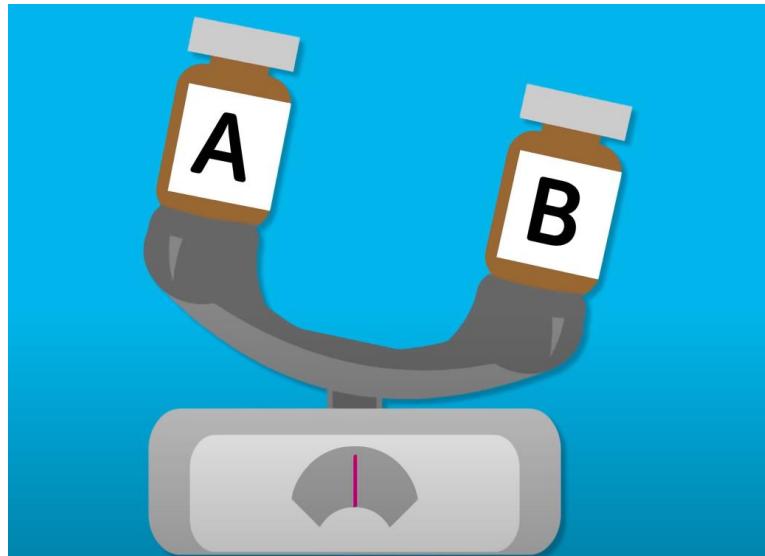


In medicine a randomized control trial, splits the population into two groups. Here, we have

- A potentially strong intervention.

Introduction (Reward on the Timeline)

Randomized Control Trials (RCT) in Medicine

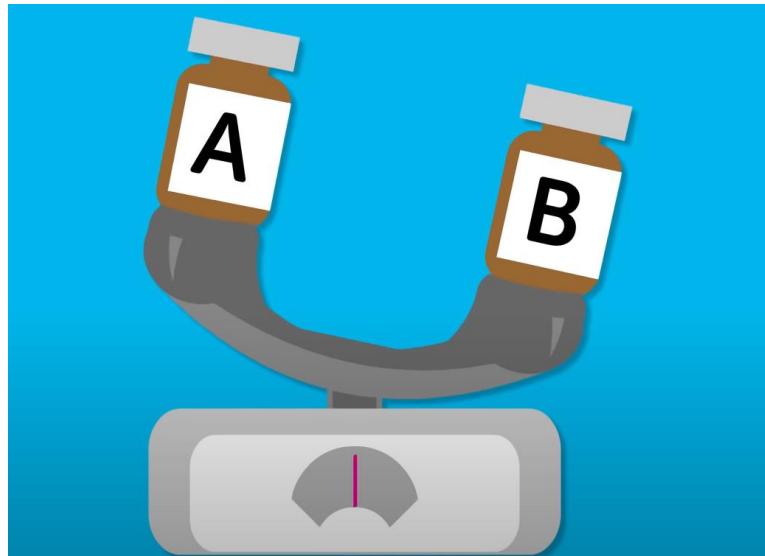


In medicine a randomized control trial, splits the population into two groups. Here, we have

- A potentially strong intervention.
- Little personalization (clear membership).

Introduction (Reward on the Timeline)

Randomized Control Trials (RCT) in Medicine

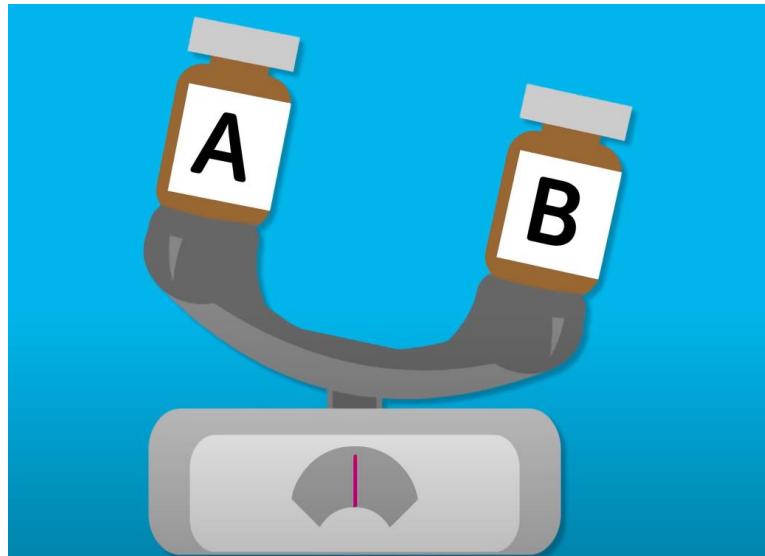


In medicine a randomized control trial, splits the population into two groups. Here, we have

- A potentially strong intervention.
- Little personalization (clear membership).
- Usually just two actions.

Introduction (Reward on the Timeline)

Randomized Control Trials (RCT) in Medicine



In medicine a randomized control trial, splits the population into two groups. Here, we have

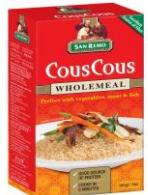
- A potentially strong intervention.
- Little personalization (clear membership).
- Usually just two actions.
- Lots of negative results!

Introduction (Reward on the Timeline)

Randomized Control Trials (RCT) in Medicine



product
viewed



Introduction (Reward on the Timeline)

Randomized Control Trials (RCT) in Medicine



product
viewed

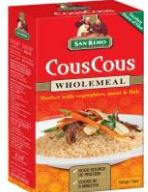


Introduction (Reward on the Timeline)

Randomized Control Trials (RCT) in Medicine



product
viewed



product
viewed



recommended
slate:

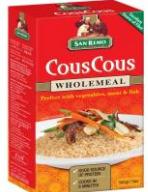


Introduction (Reward on the Timeline)

Randomized Control Trials (RCT) in Medicine



product
viewed



product
viewed



recommended
slate:



No click

Introduction (Reward on the Timeline)

Randomized Control Trials (RCT) in Medicine



product
viewed



product
viewed



recommended
slate:



product
viewed



Introduction (Reward on the Timeline)

Randomized Control Trials (RCT) in Medicine



product
viewed



product
viewed



recommended
slate:



product
viewed



recommended
slate:



Introduction (Reward on the Timeline)

Randomized Control Trials (RCT) in Medicine



product
viewed



product
viewed



recommended
slate:



product
viewed



recommended
slate:



No click

Introduction (Reward on the Timeline)

RCT in Medicine vs. A/B Tests in Recommendation

	RCT in Medicine	A/B Tests in Recommendation
Personalization segments	<ul style="list-style-type: none">The user is who they are. They are categorized by a small number of attributes e.g., age band or sex.	<ul style="list-style-type: none">The User is what they do. If they view an item - then that is who they are and what we want to personalize on.
Size of personalization segments	<ul style="list-style-type: none">Small: age bands, sex, simple medical history	<ul style="list-style-type: none">Complex – everything we know so far about the user.

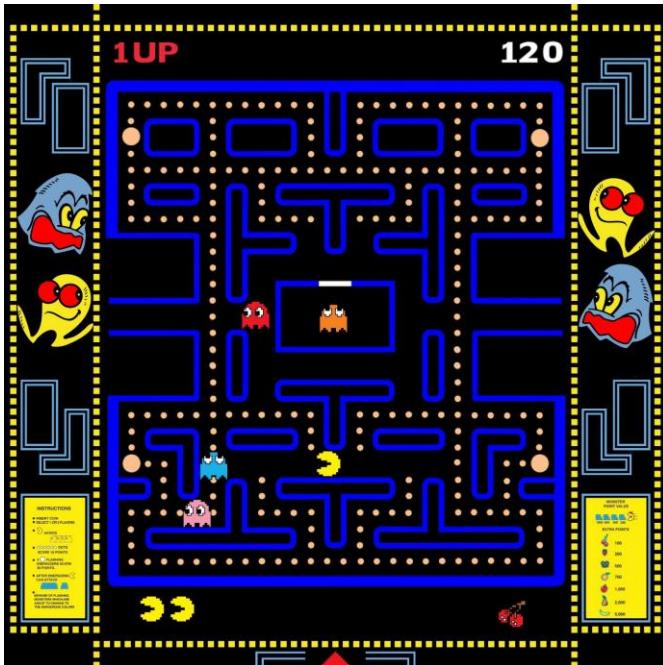
Introduction (Reward on the Timeline)

RCT in Medicine vs. A/B Tests in Recommendation

	RCT in Medicine	A/B Tests in Recommendation
Action size	<ul style="list-style-type: none">Often just intervention and placebo - (dose and dosing strategy).	<ul style="list-style-type: none">With a single item - very big. With a list (or slate) of items - the scale of the number of particles in the universe.
Effect size	<ul style="list-style-type: none">Large if it works.	<ul style="list-style-type: none">Usually relatively small.

Introduction (Reward on the Timeline)

Reinforcement Learning for Games vs Recommendation

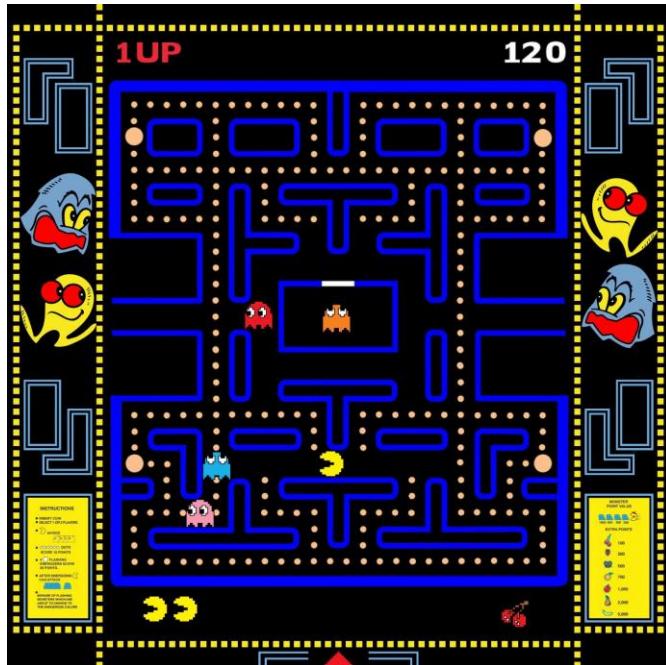


Games: In certain very successful RL settings, the problem is both similar and different to recommendation. Here, we have:

- No cost to losing many times before you win.
- Strong reward signals with little noise.
- Sequential steps of actions are critical to unlocking reward.

Introduction (Reward on the Timeline)

Reinforcement Learning for Games vs Recommendation



Recommendation: In principle it is the right formulation to move reward to the timeline. Here, we have

- High cost to performing poor actions in the wild.
- Lots of noise.
- Modest episodic features.



Introduction

Decision Theory for

Recommendation

Introduction (Decision Theory for Recommendation)

Example: Leonard Jimmie Savage in 1951



What are the **states of nature** that might occur?
→ Rain / No Rain

Introduction (Decision Theory for Recommendation)

Example: Leonard Jimmie Savage in 1951



What are the **actions** that we might take?
→ Carry umbrella / Leave umbrella

Introduction (Decision Theory for Recommendation)

Example: Leonard Jimmie Savage in 1951

Act \ State	Rain	Shine
Carry	Inconvenience and wet feet	Inconvenience and slight embarrassment
Don't carry	Miserable drenching	Bliss unalloyed



Introduction (Decision Theory for Recommendation)

Example: Leonard Jimmie Savage in 1951

Act \ State	Rain	Shine
Carry	Inconvenience and wet feet	Inconvenience and slight embarrassment
Don't carry	Miserable drenching	Bliss unalloyed



Introduction (Decision Theory for Recommendation)

Example: Leonard Jimmie Savage in 1951



Act	State	
	Rain	Shine
Carry	Inconvenience and wet feet	Inconvenience and slight embarrassment
Don't carry	Miserable drenching	Bliss unalloyed



Introduction (Decision Theory for Recommendation)

Example: Leonard Jimmie Savage in 1951



Act	State	
	Rain	Shine
Carry	Inconvenience and wet feet	Inconvenience and slight embarrassment
Don't carry	Miserable drenching	Bliss unalloyed



Introduction (Decision Theory for Recommendation)

In the Recommendation Case:

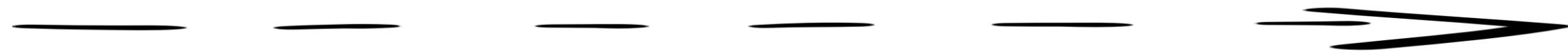
- What are the states of nature? (rain or no rain)
- What are the actions? (bring umbrella / leave umbrella)

Introduction (Decision Theory for Recommendation)

States of Nature and Decision



product
viewed



Introduction (Decision Theory for Recommendation)

States of Nature and Decision



product
viewed

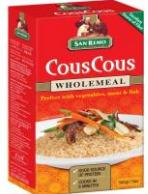


Introduction (Decision Theory for Recommendation)

States of Nature and Decision



product
viewed



product
viewed



recommended
slate:

?

Introduction (Decision Theory for Recommendation)

States of Nature and Decision



product
viewed



product
viewed



recommended
slate:

?

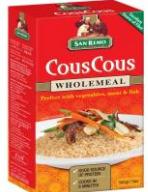
No click

Introduction (Decision Theory for Recommendation)

States of Nature and Decision



product
viewed



product
viewed



recommended
slate:

?

product
viewed

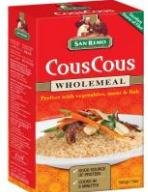


Introduction (Decision Theory for Recommendation)

States of Nature and Decision



product
viewed



product
viewed



recommended
slate:

?

product
viewed

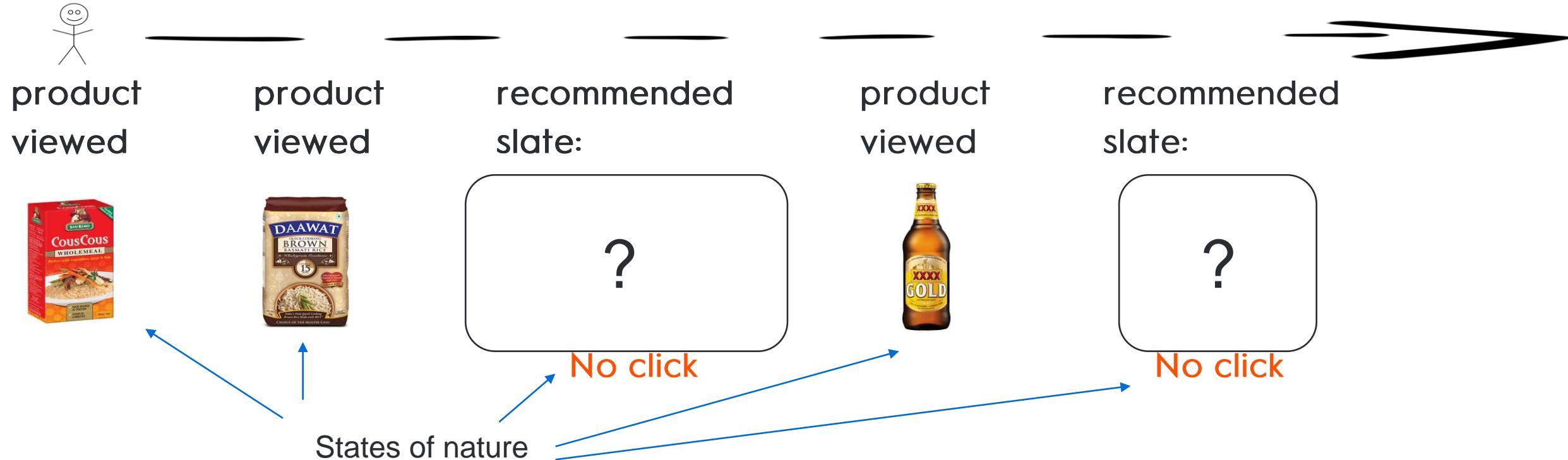


recommended
slate:

?

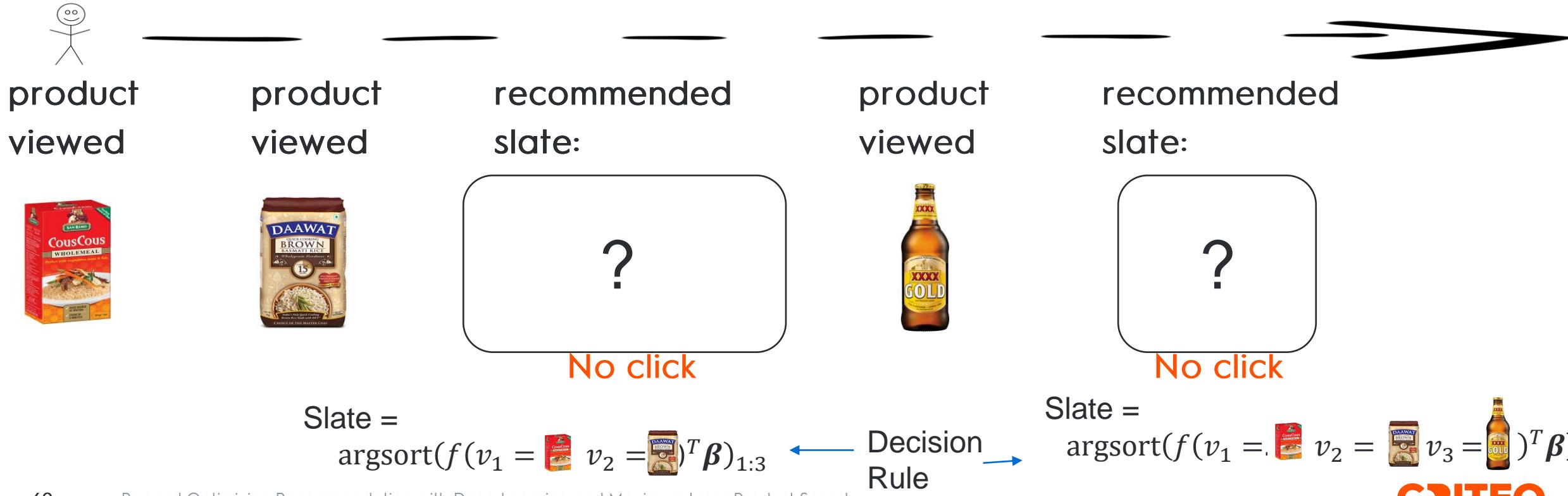
Introduction (Decision Theory for Recommendation)

States of Nature and Decision



Introduction (Decision Theory for Recommendation)

States of Nature and Decision



Introduction (Decision Theory for Recommendation)

In the Recommendation Case:

- What are the states of nature? Complex timelines of user behavior
- What are the actions? A decision rule that we can implement in production e.g., a MIPS search

The benefit of decision theory is we can restrict the complexity of the production system (decision rule) – without simplifying our understanding of the real world.

For example, we can measure the utility of two MIPS A and B by A/B test.

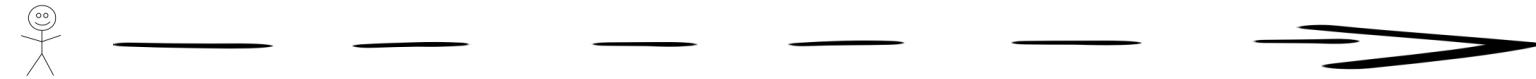
It is only possible to evaluate two candidate recommenders (driven by MIPS A and MIPS B) by understanding the probable states of nature the recommender will face. Doing this without an A/B test requires *off-policy evaluation*.

Framework

Contextual Bandits

Framework (Contextual Bandits)

Recall the full recommendation problem has the following form:



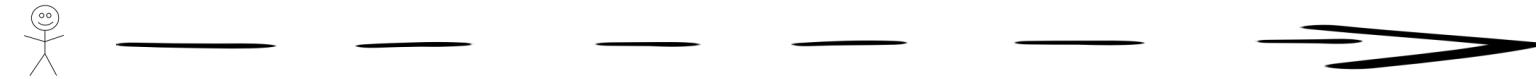
The story is:

1. We *learn something*,
2. The recommender system does something,
3. We *learn something new*,
4. ... and so on ...

Eventually we observe some reward

Framework (Contextual Bandits)

Recall the full recommendation problem has the following form:



The story is:

1. We learn something,
2. The recommender system does something,
3. We learn something new,
4. ... and so on ...

Eventually we observe some reward : this is complicated and except at A/B test time we do not know how to produce a reasonable *off policy estimate of reward*.

Framework (Contextual Bandits)

A simplified problem



The simplified story is:

1. We *learn something*,
2. The recommender system does something,
3. We *observe reward (click)*

We take slates (not users) to be conditionally independent. *Not epidemiologically correct.*

Decision Theory for Contextual Bandits

Decision theory for contextual bandits takes a particularly simplified form

x - context

R - reward

Decision rule $B(\cdot)$

Distribution on contexts $\nu(x)$

Conditional distribution on reward $R(x, a) = P[R = 1|x, a]$

The expected reward V is then given by:

$$V(B(\cdot)) = \mathbb{E}_{x \sim \nu}[R(x, B(x))]$$

Decision Theory for Contextual Bandits

Decision theory for contextual bandits takes a particularly simplified form

x - context

R - reward

Decision rule $B(\cdot)$ - we might want to constrain this to be a MIPS to be consistent with engineering constraints i.e.

$$B(x) = \text{argmax}(f_{\Xi}(x)^T \beta)$$

Distribution on contexts $\nu(x)$

Conditional distribution on reward $R(x, a) = P(R = 1|x, a)$

The expected reward V is then given by:

$$V(f_{\Xi}(\cdot), \beta) = \mathbb{E}_{x \sim \nu}[R(x, \text{argmax}(f_{\Xi}(x)^T \beta))]$$

Decision Theory for Contextual Bandits

Decision theory for contextual bandits takes a particularly simplified form

x - context

R - reward

Decision rule $B(\cdot)$ - we might further relax the decision rule to a policy

$$\pi(a|x) = \frac{\exp(f_{\Xi}(x)^T \beta_a)}{\sum_{a'} \exp(f_{\Xi}(x)^T \beta_{a'})}$$

Distribution on contexts $\nu(x)$

Conditional distribution on reward $R(x, a) = P(R = 1|x, a)$

The expected reward V is then given by:

$$V(f_{\Xi}(\cdot), \beta) = \mathbb{E}_{x \sim \nu} \mathbb{E}_{a \sim \pi(\cdot|x)} [R(x, a)]$$

Decision Theory for Contextual Bandits

Decision theory for contextual bandits takes a particularly simplified form

x - context

R - reward

Decision rule $B(\cdot)$ - we might further relax the decision rule to a policy

$$\pi(a|x) = \frac{\exp(f_{\Xi}(x)^T \beta_a)}{\sum_{a'} \exp(f_{\Xi}(x)^T \beta_{a'})}$$

Distribution on contexts $\nu(x)$

Conditional distribution on reward $R(x, a) = P(R = 1|x, a)$

The expected reward V is then given by:

$$V(f_{\Xi}(\cdot), \beta) = V(\pi) = \mathbb{E}_{x \sim \nu} \mathbb{E}_{a \sim \pi(\cdot|x)} [R(x, a)]$$

Our fundamental problem is to find the optimal MIPS i.e., $f_{\Xi}(\cdot), \beta$. In order to evaluate this, we need to understand how the world works i.e., $\nu(x)$ and particularly $P(R = 1|x, a) = R(x, a)$.

$R(x, a)$ is particularly important because the best a^* for a given x is given by:

$$a^* = \operatorname{argmax}_a R(a, x)$$

Decision Theory for Contextual Bandits

Decision theory for contextual bandits takes a particularly simplified form

x - context

R - reward

Decision rule $B(\cdot)$ - we might further relax the decision rule to a policy

$$\pi(a|x) = \frac{\exp(f_{\Xi}(x)^T \beta_a)}{\sum_{a'} \exp(f_{\Xi}(x)^T \beta_{a'})}$$

Distribution on contexts $\nu(x)$

Conditional distribution on reward $R(x, a) = P(R = 1|x, a)$

The expected reward V is then given by:

$$V(f_{\Xi}(\cdot), \beta) = V(\pi) = \mathbb{E}_{x \sim \nu} \mathbb{E}_{a \sim \pi(\cdot|x)} [R(x, a)]$$

Key reward optimizing recommendation approaches focus on:

- Off policy estimation of $R(x, a)$.
- Given access to an estimate of $R(x, a)$ how do we get the decision rule $f_{\Xi}(\cdot), \beta$.

Successful reward optimizing recommendation requires both.

Framework Off-Policy Reward Estimation

Framework (Off-Policy Reward Estimation)

Two Families of Approaches

- The goal of a reward optimizing recommender system is to find a policy π that maximizes the value function

$$V(\pi) = \mathbb{E}_{x \sim \nu} \mathbb{E}_{a \sim \pi(\cdot|x)} [R(x, a)]$$

- We never have access to the true value function $V(\pi)$, thus the first step is to find an estimate $\hat{V}_n(\pi)$ of $V(\pi)$. There are two families of approaches for this purpose, Direct Methods and Inverse Propensity Scoring.

Framework (Off-Policy Reward Estimation)

Two Families of Approaches

Direct Methods (DM): Find a model-based estimator $\hat{R}^{DM}(x, a)$ of the expected reward $R(x, a)$. The estimator of the value function follows as

$$\hat{V}_n^{DM}(\pi) = \frac{1}{n} \sum_{i=1}^n \sum_a \hat{R}^{DM}(x_i, a) \pi(a | x_i)$$

Inverse Propensity Scoring (IPS): Design unbiased estimators of $V(\pi)$ by re-weighting samples using the importance sampling trick

$$\hat{V}_n^{IPS}(\pi) = \frac{1}{n} \sum_{i=1}^n R_i \frac{\pi(a_i | x_i)}{\pi_0(a_i | x_i)}$$

Framework (Off-Policy Reward Estimation)

Two Families of Approaches

Direct Methods (DM): $\hat{V}_n^{DM}(\pi) = \frac{1}{n} \sum_{i=1}^n \sum_a \hat{R}^{DM}(x_i, a) \pi(a | x_i)$

Inverse Propensity Scoring (IPS): $\hat{V}_n^{IPS}(\pi) = \frac{1}{n} \sum_{i=1}^n R_i \frac{\pi(a_i | x_i)}{\pi_0(a_i | x_i)}$

→ IPS can be seen as a DM where :

$$\hat{R}^{IPS}(x, a) = \begin{cases} \frac{1}{|\mathcal{I}_n(x, a)|} \sum_{i \in \mathcal{I}_n(x, a)} \frac{R_i}{\pi_0(a_i | x_i)} & \text{if } |\mathcal{I}_n(x, a)| \geq 1 \\ 0 & \text{otherwise} \end{cases}$$

where $\mathcal{I}_n(x, a) = \{i \in [n], x_i = x \text{ and } a_i = a\}$

Framework (Off-Policy Reward Estimation)

Two Families of Approaches

Direct Methods (DM): $\hat{V}_n^{DM}(\pi) = \frac{1}{n} \sum_{i=1}^n \sum_a \hat{R}^{DM}(x_i, a) \pi(a | x_i)$

Inverse Propensity Scoring (IPS): $\hat{V}_n^{IPS}(\pi) = \frac{1}{n} \sum_{i=1}^n R_i \frac{\pi(a_i | x_i)}{\pi_0(a_i | x_i)}$

- $\hat{R}(x_i, a)$ will denote an estimator of the expected reward $R(x, a)$ using either DM or IPS.
- Similarly, we refer by $\hat{V}_n(\pi)$ an estimator of the value function using either $\hat{V}_n^{DM}(\pi)$ or $\hat{V}_n^{IPS}(\pi)$.

$$\hat{V}_n(\pi) = \frac{1}{n} \sum_{i=1}^n \sum_a \hat{R}(x_i, a) \pi(a | x_i)$$



Single Item Recommendation

IPS Method

Single Item Recommendation (IPS Method)

IPS for Single Items

- $\hat{V}_n^{IPS}(\pi)$ is provably unbiased under mild assumptions.
- However, the variance of $\hat{V}_N^{IPS}(\pi)$ grows linearly with the importance weight $\frac{\pi(a_i|x_i)}{\pi_0(a_i|x_i)}$.
→ the variance explodes when π drifts away from π_0 .
- Example: The logging policy π_0 is uniform. In this case, the variance grows with $P\pi(a|x)$, where P is the catalog size (massive for large catalogs).

Single Item Recommendation (IPS Method)

IPS for Single Items

- Existing solutions of the variance problem of $\hat{V}_n^{IPS}(\pi)$ include:
 - Clipping the importance weights: replacing $\frac{\pi(a|x)}{\pi_0(a|x)}$ by $\min\left(\frac{\pi(a|x)}{\pi_0(a|x)}, M\right)$.
 - Self normalization of the importance weights.
 - Combining DM and IPS using doubly robust estimation (DR).
- Introduce a **bias-variance** tradeoff but still difficult when **P is large**.

Dudík, Miroslav, et al. 2012.

Bottou, Léon, et al. 2013

Swaminathan, Adith, and Thorsten Joachims. 2015



Single Item Recommendation Direct Method

Single Item Recommendation (Direct Method)

Logistic Regression Model

$$R(\mathbf{x}, a) = P(R = 1 | \mathbf{x}, a) = \sigma(g_{\Gamma}(\mathbf{x})^T \boldsymbol{\Psi}_a)$$

- Can be biased but doesn't suffer from a variance problem.
- **Reduce the bias** by increasing the embeddings dimension D.

BLOB – uses priors to accelerate learning the reward in large action spaces.

Because of the parameterization of the reward model, we automatically get a MIPS:

$$a^* = \operatorname{argmax}_a (f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_a)$$

with $f_{\Xi}(\cdot) \leftarrow g_{\Gamma}(\cdot)$ and $\boldsymbol{\beta} \leftarrow \boldsymbol{\Psi}$

it can result in larger than
necessary embeddings

Single Item Recommendation (Direct Method)

Logistic Regression Model

$$R(\mathbf{x}, a) = P(R = 1 | \mathbf{x}, a) = \sigma(g_{\Gamma}(\mathbf{x})^T \boldsymbol{\Psi}_a)$$

- Can be biased but doesn't suffer from a variance problem.
- **Reduce the bias** by increasing the embeddings dimension D.

BLOB – uses priors to accelerate learning the reward in large action spaces.

Because of the parameterization of the reward model, we automatically get a MIPS:

$$a^* = \operatorname{argmax}_a (f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_a)$$

with $f_{\Xi}(\cdot) \leftarrow g_{\Gamma}(\cdot)$ and $\boldsymbol{\beta} \leftarrow \boldsymbol{\Psi}$

We will discuss alternative ways to get $f_{\Xi}(\cdot), \boldsymbol{\beta}$

Single Item Recommendation (Direct Method)

Logistic Regression Model

We can split $x = [y, z]$ where

y - User Engagement, and z - User interest

$$R(x, a) = P(R = 1 | x, a) = \sigma(y^T \phi + g_\Gamma(z)^T \Psi_a)$$

$$a^* = \operatorname{argmax}_a (f_\Xi(z)^T \beta_a)$$

with $f_\Xi(\cdot) \leftarrow g_\Gamma(\cdot)$ and $\beta \leftarrow \Psi$

Single Item Recommendation (Direct Method)

Logistic Regression Model

We can split $x = [y, z]$ where

y - User Engagement, and z - User interest

$$R(x, a) = P(R = 1 | x, a) = \sigma(y^T \phi + g_\Gamma(z)^T \Psi_a)$$

$$a^* = \operatorname{argmax}_a (f_\Xi(z)^T \beta_a)$$

Nuisance parameters

with $f_\Xi(\cdot) \leftarrow g_\Gamma(\cdot)$ and $\beta \leftarrow \Psi$

Single Item Recommendation (Direct Method)

Remarks



product
viewed



Single Item Recommendation (Direct Method)

Remarks



product
viewed

product
viewed



Single Item Recommendation (Direct Method)

Remarks



product
viewed



product
viewed



recommended
item:



Single Item Recommendation (Direct Method)

Remarks



product
viewed



product
viewed



recommended
 slate item:



Click/No click

Single Item Recommendation (Direct Method)

Remarks



product
viewed



product
viewed



recommended
 slate item:



Click/No click

Recommend =

$$\text{argmax}(f(v_1 = \text{rice}, v_2 = \text{rice}))^T \boldsymbol{\beta}$$

Single Item Recommendation (Direct Method)

Remarks



product
viewed



product
viewed



recommended
state item:



Click/No click

Recommend =

$$\text{argmax}(f(v_1 = \text{rice}, v_2 = \text{rice}))^T \boldsymbol{\beta}$$

How do we build the reward model?



We can also use features like:

- How active is the user?
- How big/attractive is the recommendation?

Single Item Recommendation (Direct Method)

Remarks



product
viewed



product
viewed



recommended
state item:



Click/No click

Recommend =

$$\text{argmax}(f(v_1 = \text{rice}, v_2 = \text{rice})^T \boldsymbol{\beta})$$

How do we build the reward model?



We can also use features like:

- How active is the user?
- How big/attractive is the recommendation?

This (often)
dominates
in practice

Single Item Recommendation (Direct Method)

Remarks



product
viewed



product
viewed



recommended
state item:



Click/No click

Recommend =

$$\text{argmax}(f(v_1 = \text{[image of CousCous]}, v_2 = \text{[image of Brown Rice]}))^T \boldsymbol{\beta})$$

How do we build the reward model?



We can also use features like:

- How active is the user?
- How big/attractive is the recommendation?

We can predict
well using only
these features



Single Item Recommendation Decision Rule Optimization

Single Item Recommendation (Decision Rule Optimization)

Why Optimize $\pi(a | x)$ when $R(x, a)$ is a Model?

If our reward model has the following parametric form:

$$\hat{R}(x, a) = P(R = 1 | x, a) = \sigma(g_\Gamma(x)^T \Psi_a)$$

Then we already have an optimal MIPS: $f_\Xi(x)^T \beta$, by setting

$$f_\Xi(x) \leftarrow g_\Gamma(x), \quad \beta \leftarrow \Psi$$

So why would we optimize: $\hat{V}_n(\pi) = \sum_{i=1}^n \sum_a \hat{R}(x_i, a) \pi(a | x_i)$ to obtain a MIPS if we already have one?

→ $\hat{R}(x, a) = P(R = 1 | x, a) = \sigma(g_\Gamma(x)^T \Psi_a)$ is sufficient but not necessary.

→ A few reasons we might want not to use it in the decision rule.

Single Item Recommendation (Decision Rule Optimization)

Why Optimize $\pi(a | x)$ When $R(x, a)$ is a Model?

→ A few reasons we might want not to use it in the decision rule

- $\sigma(y^T \phi + f_{\Xi}(x)^T \beta_a)$ is a calibrated reward that correctly orders all actions. We only need the best action.
- Imposing the above can mean the embedding dimension D for $f_{\Xi}(x)^T \beta$ might be larger than we need.
- We might want the dimension of β lower than D for engineering constraints.
- We might prefer a more general model $P(R = 1 | a, x) = \sigma(h_{\Phi}(x, a))$.

Single Item Recommendation (Decision Rule Optimization)

Optimizing the decision rule via a policy

The objective is:

$$\hat{V}_n(\pi_{\Xi, \beta}) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{a \sim \pi_{\Xi, \beta}(\cdot | \mathbf{x})} [\hat{R}(a, \mathbf{x})]$$

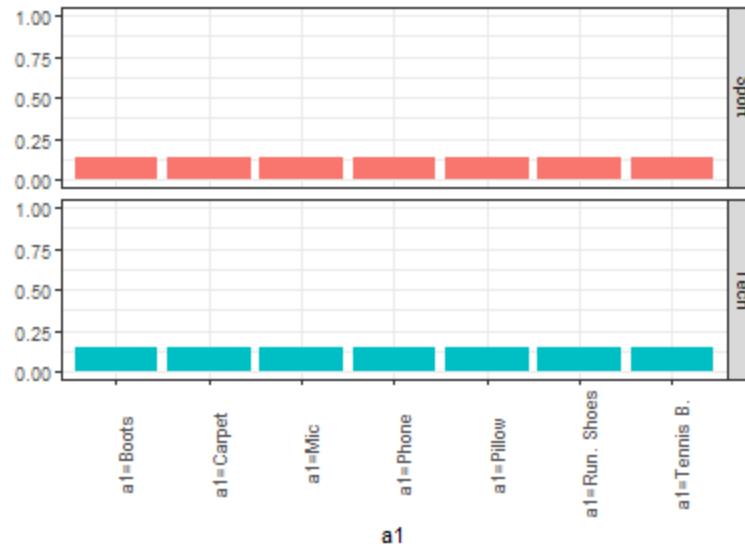
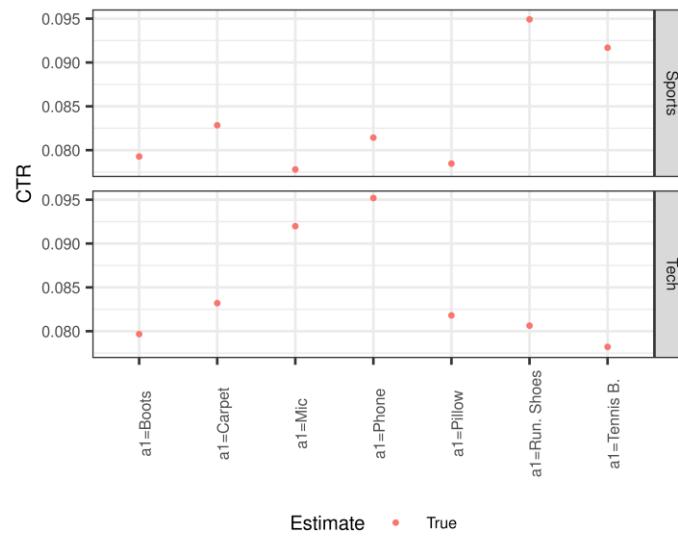
where :

$$\pi_{\Xi, \beta}(a | \mathbf{x}) = \frac{\exp(f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_a)}{\sum_i \exp(f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_i)}$$

If we want $\pi(a | \mathbf{x})$ to be restricted to producing a MIPS

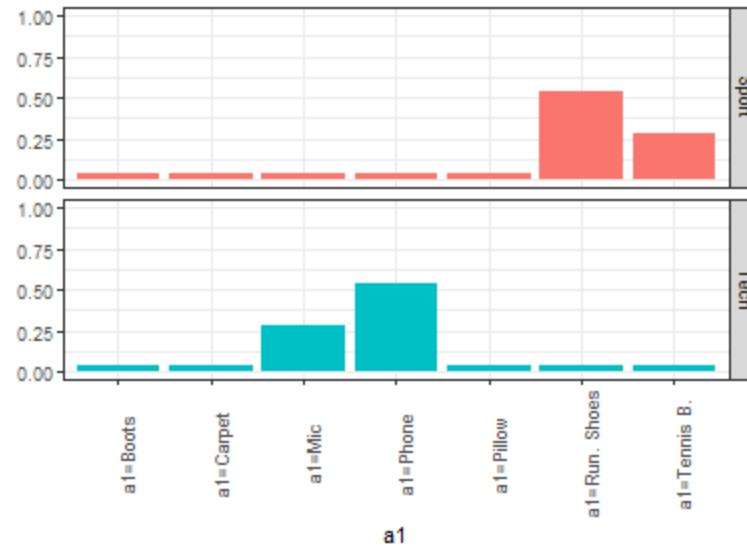
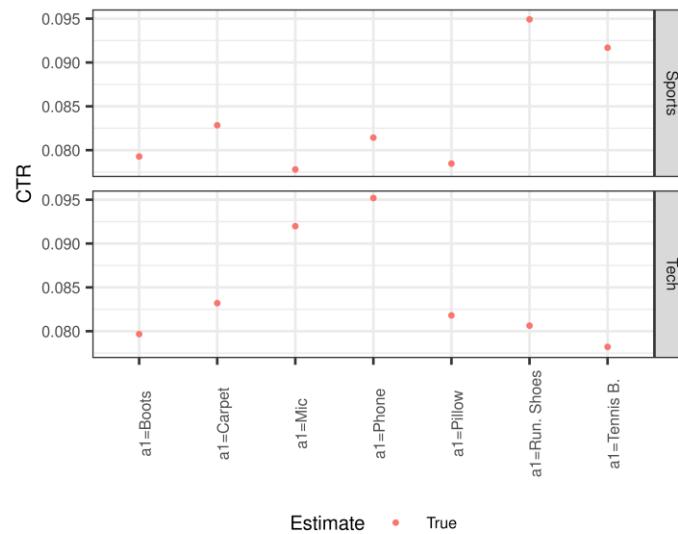
Single Item Recommendation (Decision Rule Optimization)

Optimizing the decision rule via a policy



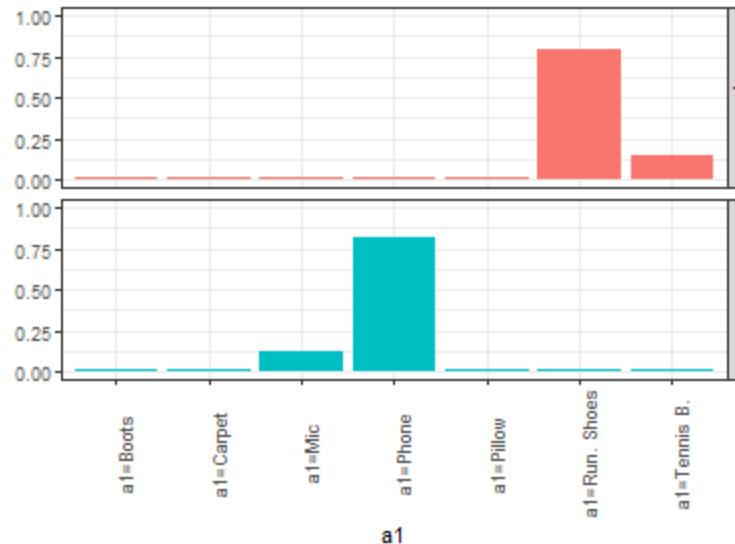
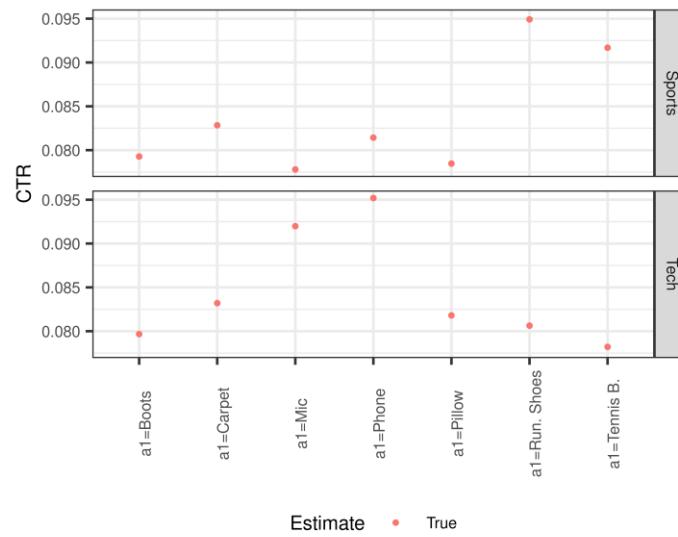
Single Item Recommendation (Decision Rule Optimization)

Optimizing the decision rule via a policy



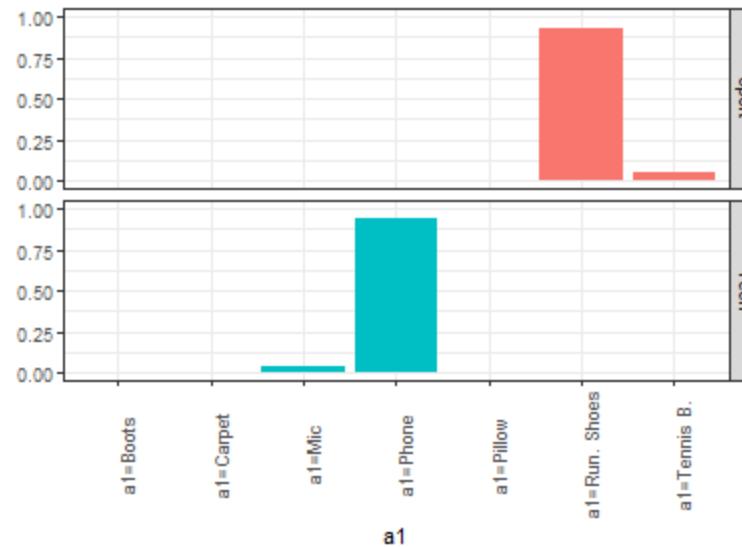
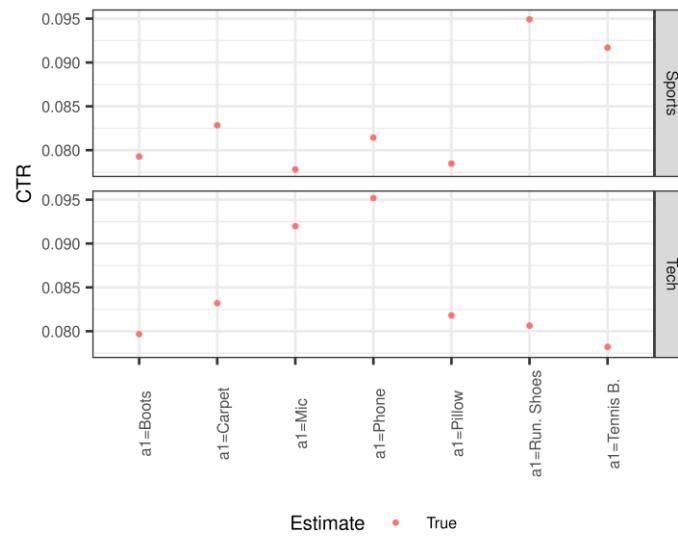
Single Item Recommendation (Decision Rule Optimization)

Optimizing the decision rule via a policy



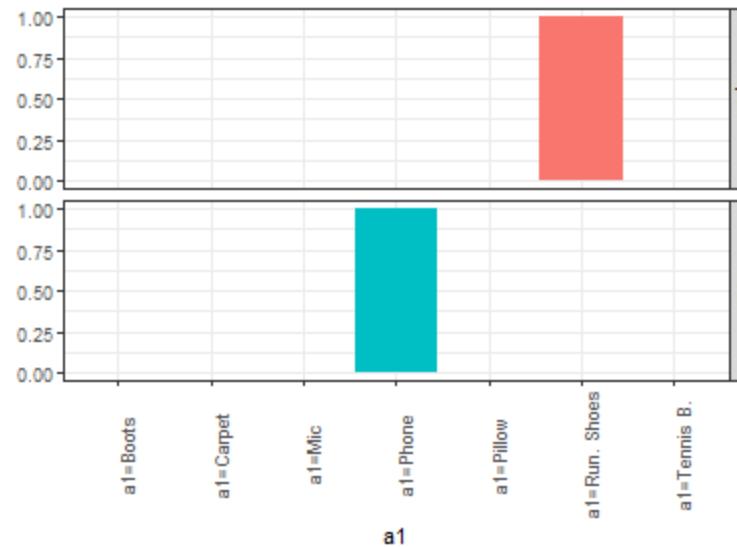
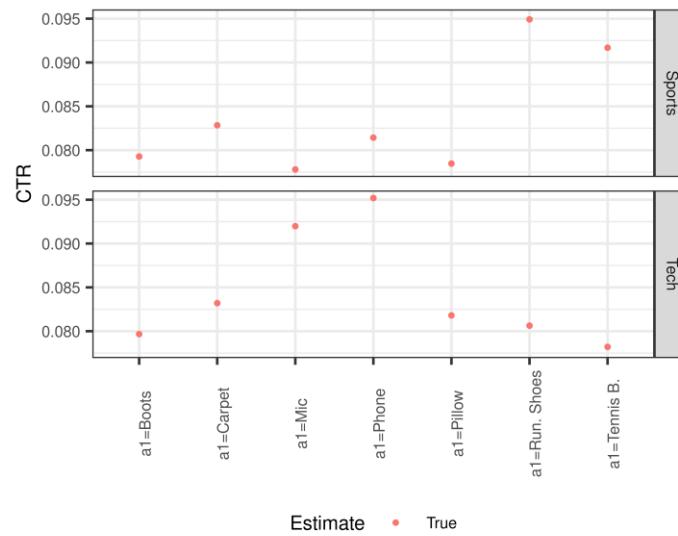
Single Item Recommendation (Decision Rule Optimization)

Optimizing the decision rule via a policy



Single Item Recommendation (Decision Rule Optimization)

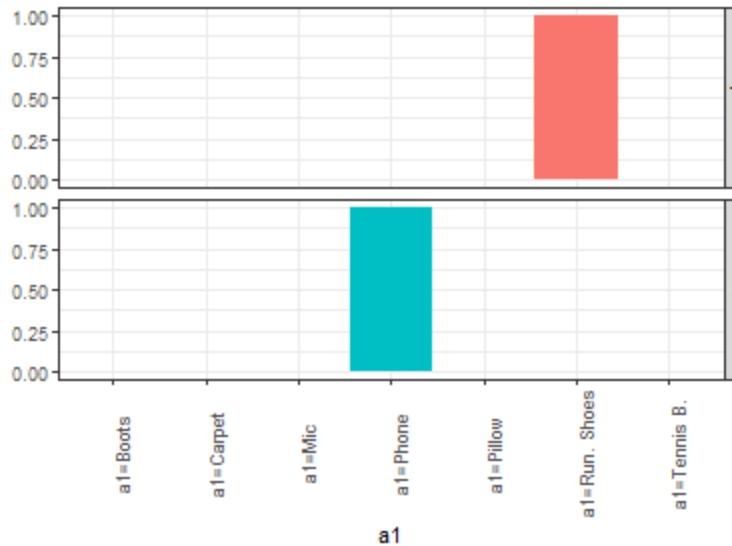
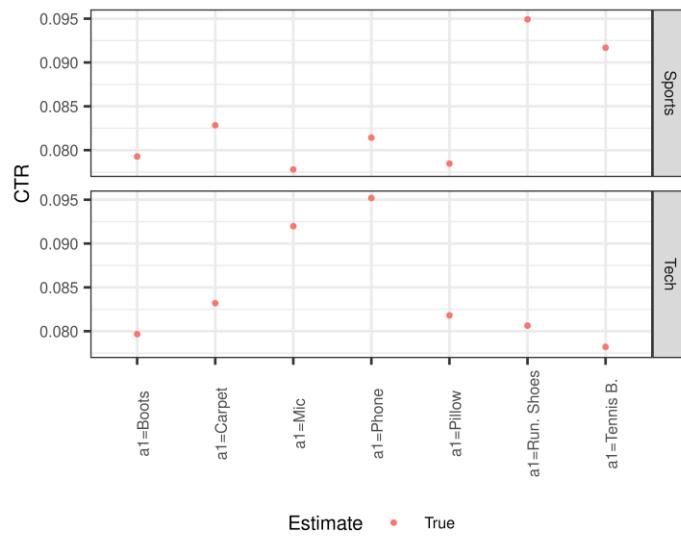
Optimizing the decision rule via a policy



At convergence $\pi(a|x)$ is deterministic.

Single Item Recommendation (Decision Rule Optimization)

Optimizing the decision rule via a policy



At convergence $\pi(a|x)$ is deterministic.

If the reward estimate is wrong, it will (possibly) select the wrong best action.

Single Item Recommendation (Decision Rule Optimization)

Optimizing the MIPS with Policy Learning

$$\hat{V}_{\Xi, \beta}(\mathbf{x}) = \mathbb{E}_{a \sim \pi_{\Xi, \beta}(\cdot | \mathbf{x})} [\hat{R}(a, \mathbf{x})] \text{ where } \pi_{\Xi, \beta}(a | \mathbf{x}) = \frac{\exp(f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_a)}{\sum_i \exp(f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_i)}$$

Exact Gradients:

$$\nabla_{\Xi, \beta} \hat{V}_{\Xi, \beta}(\mathbf{x}) = \nabla_{\Xi, \beta} \mathbb{E}_{a \sim \pi_{\Xi, \beta}(\cdot | \mathbf{x})} [\hat{R}(a, \mathbf{x})] = \mathbb{E}_{a \sim \pi_{\Xi, \beta}(\cdot | \mathbf{x})} [\hat{R}(a, \mathbf{x}) \nabla_{\Xi, \beta} \log \pi(a | \mathbf{x})]$$

Swaminathan, Adith, and Thorsten Joachims. 2015.

Joachims, Thorsten, Adith Swaminathan, and Maarten De Rijke. 2018.

Single Item Recommendation (Decision Rule Optimization)

Optimizing the MIPS with REINFORCE

$$\hat{V}_{\Xi, \beta}(\mathbf{x}) = \mathbb{E}_{a \sim \pi_{\Xi, \beta}(\cdot | \mathbf{x})} [\hat{R}(a, \mathbf{x})] \text{ where } \pi_{\Xi, \beta}(a | \mathbf{x}) = \frac{\exp(f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_a)}{\sum_i \exp(f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_i)}$$

Approximate Gradients:

$$a' \sim \pi_{\Xi, \beta}(\cdot | \mathbf{x})$$

$$\nabla_{\Xi, \beta} \hat{V}_{\Xi, \beta}(\mathbf{x}) = \hat{R}(\mathbf{x}, a') \nabla_{\Xi, \beta} \log \pi(a' | \mathbf{x})$$

On its own we don't get a speed up for using this algorithm.

Chen, et al. 2018.



Single Item Recommendation

SNIPS REINFORCE

Single Item Recommendation (SNIPS REINFORCE)

Optimizing the Objective

Our objective is:

$$\nabla_{\Xi, \beta} \hat{V}_{\Xi, \beta}(x) = \nabla_{\Xi, \beta} \mathbb{E}_{a \sim \pi_{\Xi, \beta}(\cdot | x)} [\hat{R}(x, a)] = \mathbb{E}_{a \sim \pi_{\Xi, \beta}(\cdot | x)} [\hat{R}(x, a) \nabla_{\Xi, \beta} \log \pi_{\Xi, \beta}(a | x)]$$

There are numerous problems with this gradient :

- $\log \pi_{\Xi, \beta}(a | x)$ involves the computation of the normalizing constant $\rightarrow O(P)$
- If we want to compute the exact expectation $\rightarrow O(P)$ (Infeasible)
- If we want to approximate the gradient by sampling $\rightarrow O(P)$

Everything scales **at least linearly** in P which can be very costly in large catalog sizes.

Single Item Recommendation (SNIPS REINFORCE)

Optimizing the Objective is not Maximum Likelihood

$$\log P(a|x, \beta, \Xi) = f_{\Xi}(x)^T \beta_a - \log \sum_{a'} \exp(f_{\Xi}(x)^T \beta_{a'})$$

$$\nabla_{\beta, \Xi} \log P(a|x, \beta, \Xi) = \nabla_{\beta, \Xi} f_{\Xi}(x)^T \beta_a - \mathbb{E}_{P(a'|x, \beta, \Xi)} [\nabla_{\beta, \Xi} f_{\Xi}(x)^T \beta_{a'}]$$

Lots of work on how to solve this problem! e.g.,

Bengio, Yoshua, and Jean-Sébastien Senécal. 2003.

Ruiz, Francisco, et al. 2018.

Rawat, Ankit Singh, et al. 2019.

Single Item Recommendation (SNIPS REINFORCE)

Optimizing the Objective

Our objective is:

$$\hat{V}_{\Xi, \beta}(\mathbf{x}) = \mathbb{E}_{a \sim \pi_{\Xi, \beta}(\cdot | \mathbf{x})} [\hat{R}(a, \mathbf{x})] \quad \text{where } \pi_{\Xi, \beta}(a | \mathbf{x}) = \frac{\exp(f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_a)}{\sum_i \exp(f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_i)}$$
$$\nabla_{\Xi, \beta} \hat{V}_{\Xi, \beta}(\mathbf{x}) = \mathbb{E}_{a \sim \pi_{\Xi, \beta}(\cdot | \mathbf{x})} [\hat{R}(\mathbf{x}, a) \nabla_{\Xi, \beta} \log \pi_{\Xi, \beta}(a | \mathbf{x})]$$

→ We propose a new gradient :

$$\nabla_{\Xi, \beta} \hat{V}_{\Xi, \beta}(\mathbf{x}) = Cov_{a \sim \pi_{\Xi, \beta}(\cdot | \mathbf{x})} [\hat{R}(\mathbf{x}, a), \nabla_{\Xi, \beta} f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_a]$$

with $Cov[Y, \mathbf{Z}] = \mathbb{E}[(Y - \mathbb{E}[Y])(\mathbf{Z} - \mathbb{E}[\mathbf{Z}])]$

Single Item Recommendation (SNIPS REINFORCE)

Self Normalized Importance Sampling

- Our new gradient gets rid of the normalizing constant, is that sufficient?

$$\nabla_{\Xi, \beta} \hat{V}_{\Xi, \beta}(\mathbf{x}) = Cov_{a \sim \pi_{\Xi, \beta}(\cdot | \mathbf{x})} [\hat{R}(\mathbf{x}, a), \nabla_{\Xi, \beta} f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_a]$$

NO! We still need to avoid sampling from $\pi_{\Xi, \beta}(\cdot | \mathbf{x})$ to approximate the gradient.

Solution : Self Normalized Importance Sampling!

Single Item Recommendation (SNIPS REINFORCE)

Self Normalized Importance Sampling

If we have a proposal Q that is fast to sample from, we can approximate the covariance gradient $Cov[Y, Z]$ by $\sum_{i=1}^S w_i(Y_i - \bar{Y})(Z_i - \bar{Z})$ with samples a_1, \dots, a_S from Q with :

$$\tilde{\pi}(a|x) = \exp(f_{\Xi}(x)^T \beta_a), \quad Y_i = Y(a_i), \quad Z_i = Z(a_i)$$

$$\tilde{w}_i = \frac{\tilde{\pi}(a_i|x)}{Q(a_i|x)}, \quad w_i = \frac{\tilde{w}_i}{\sum \tilde{w}_j}, \quad \bar{Y} = \sum_{i=1}^N w_i Y_i, \quad \bar{Z} = \sum_{i=1}^N w_i Z_i$$

→ How do we choose Q?

Single Item Recommendation (SNIPS REINFORCE)

Proposal Selection

We need an additional Assumption :

$$\pi_{\Xi, \beta}(a|x) = \frac{\exp(f_{\Xi}(x)^T \beta)}{\sum_{a'} \exp(f_{\Xi}(x)^T \beta_{a'})}$$

Assume β is **fixed** and we only optimize $f_{\Xi}(x)$. This structure allows fast search in $O(\log P)$.

MIPS give us $\alpha_K(x) = \text{argsort}(f_{\Xi}(x)^T \beta)_{1:K}$ the top-K (approximate) items of the policy

We define :

$$Q_{\epsilon, K}(a|x) = \begin{cases} \frac{\epsilon}{P} + (1 - \epsilon)\kappa(a|x) & \text{if } a \in \alpha_K(x) \\ \frac{\epsilon}{P} & \text{else} \end{cases} \quad \text{with } \kappa(a|x) = \frac{\exp(f_{\Xi}(x)^T \beta_a)}{\sum_{a' \in \alpha_K(x)} \exp(f_{\Xi}(x)^T \beta_{a'})}$$

Single Item Recommendation (SNIPS REINFORCE)

Experiments

We test the method on a session completion task : X observed session, Y complementary items.

This can be cast into an offline bandit framework where our policy π_θ takes the observed part X as a context, recommends an item a and receives the binary reward $\hat{r}(a, X) = 1[a \in Y]$.

We use X to build item embeddings β of dimension $L \ll P$.

Define x_{emb} as the item embeddings average in session X.

$$\pi(a|x) = \frac{\exp(f_\Xi(x)^T \beta_a)}{\sum_{a'} \exp(f_\Xi(x)^T \beta_{a'})}$$

	Catalog size	Number of users
Twitch	750K	500K
GoodReads	1.23M	300K

Single Item Recommendation (SNIPS REINFORCE)

Experiments

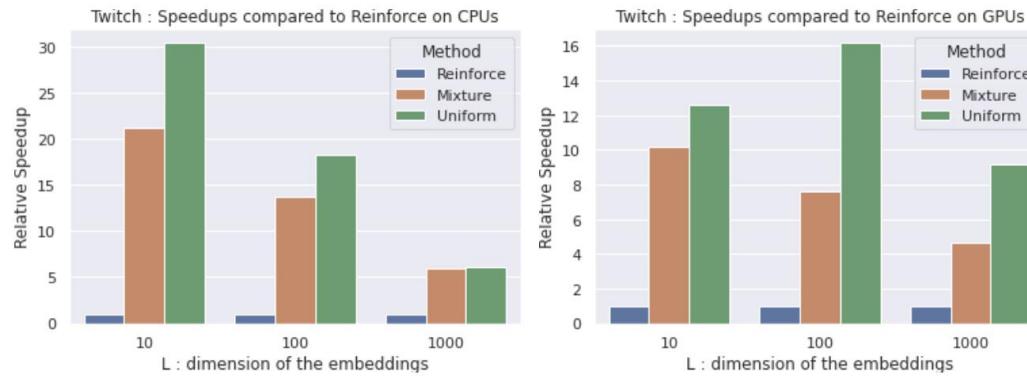
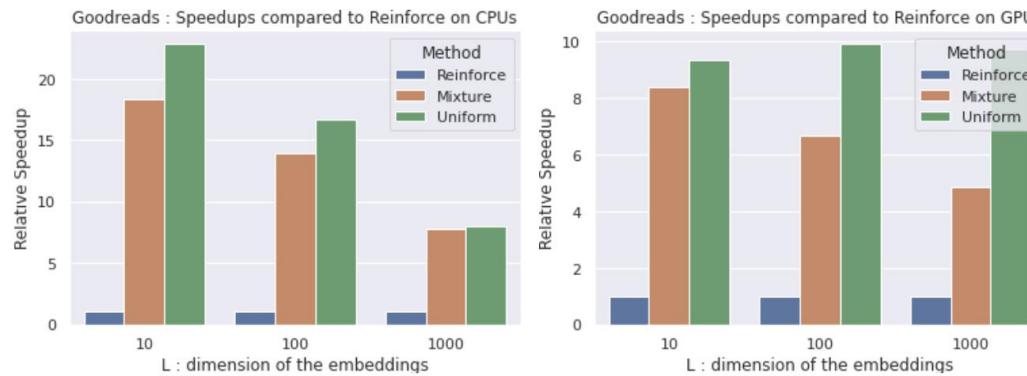


Fig. 1. The Twitch Dataset : The speedup of the proposed algorithms on both CPU and GPU devices



Single Item Recommendation (SNIPS REINFORCE)

Experiments

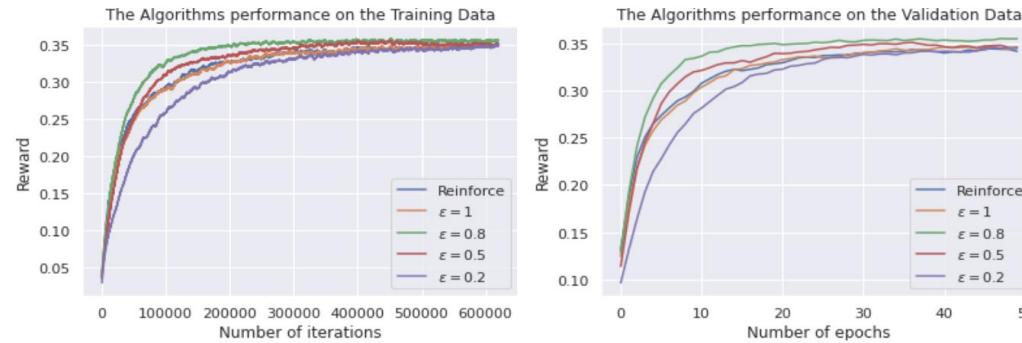


Fig. 3. The Twitch Dataset : ϵ effect on the policy training and validation

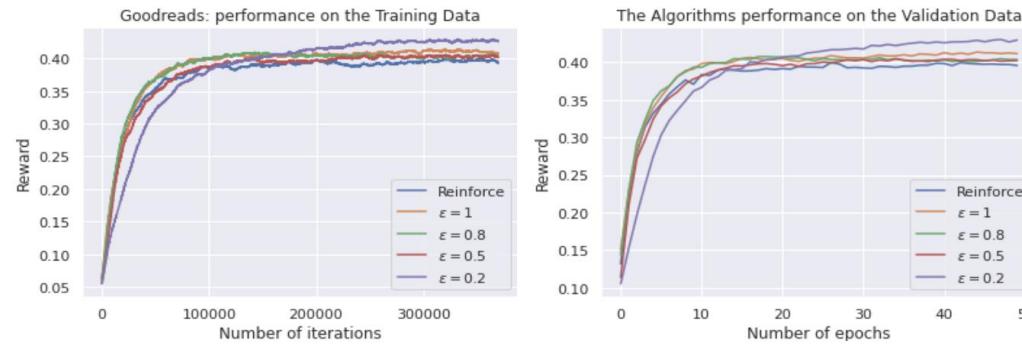


Fig. 4. The Goodreads Dataset : ϵ effect on the policy training and validation

Single Item Recommendation (SNIPS REINFORCE)

Experiments

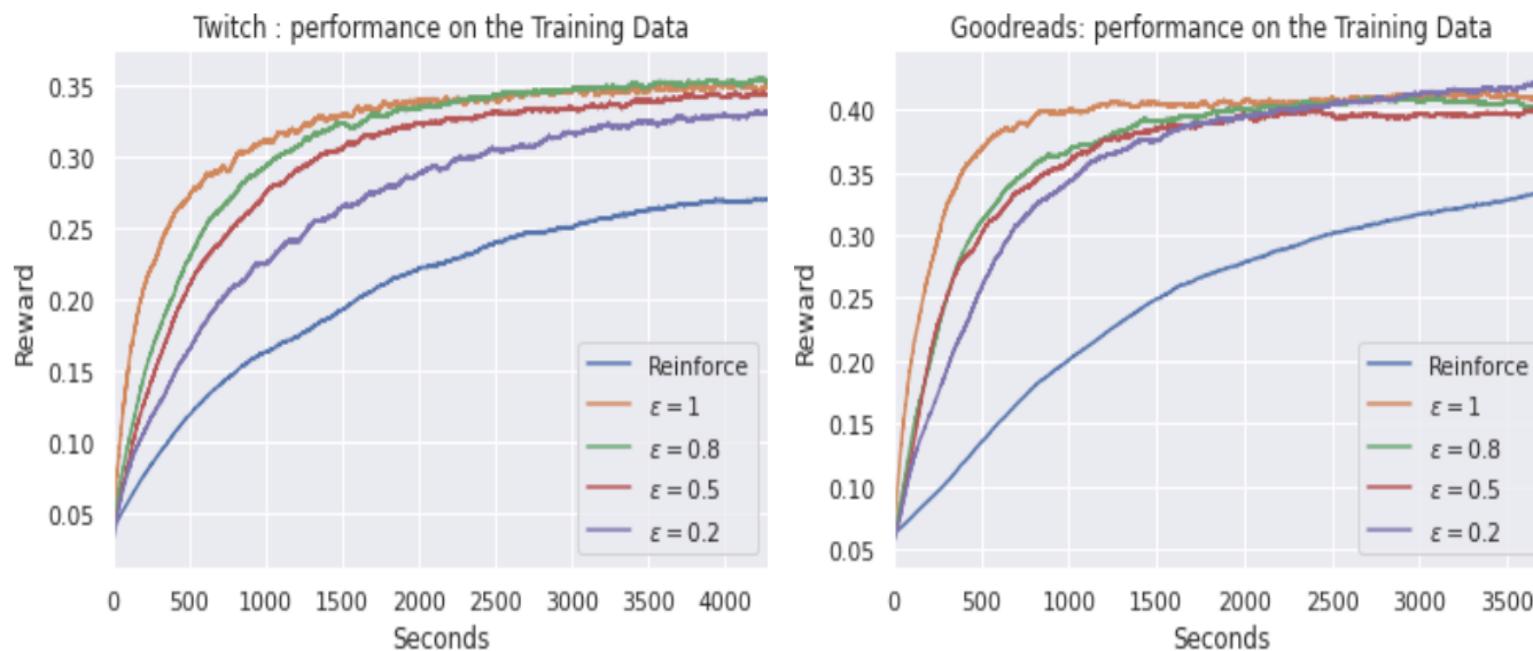


Fig. 5. Performance plots given time : ϵ 's effect on the training

Sakhi et al. 2022

SNIPS Reinforce (Practice)

Python Tutorial 2

https://colab.research.google.com/github/otmhi/Reward-Optimizing-Reco/blob/main/Reward_Optimizing_Slate_Recommendation_with_DL_and_MIPS_Part_2.ipynb



Slate

Recommendation

Motivation

Slate Recommendation (Motivation)

Recommender Systems in Real-World



product
viewed



Slate Recommendation (Motivation)

Recommender Systems in Real-World



product
viewed



Slate Recommendation (Motivation)

Recommender Systems in Real-World



product
viewed

product
viewed

recommended
slate:



Slate Recommendation (Motivation)

Recommender Systems in Real-World



product
viewed



product
viewed



recommended
slate:



No click

No click

click



Slate

Recommendation

Framework

Slate Recommendation (Framework)

From Single Items to Slates

- The recommender system observes a context $x \sim \nu$.
 - The recommender system selects an ordered list (or slate) of K items a_1, a_2, \dots, a_K .
 - The recommender system receives a reward $R \sim P_R(\cdot | x, a_1, \dots, a_K)$ and a list of ranks $r_1, r_2, \dots, r_K \sim P_r(\cdot | x, s)$.
- Similarly, to single item recommendation, the recommender system selects a slate of K items a_1, a_2, \dots, a_K following a decision rule $B(\cdot)$ that boils down to a MIPS task.

Slate Recommendation (Framework)

From Single Items to Slates

- The space of possible actions (i.e., slates) is huge compared to single item recommendation (P^K).
- Estimating the reward is more difficult, the reward signal is spread over a huge combination of context-slate pairs.
- Additional biases related to the layout of the slates (e.g., position bias).
- Additional factors that might have a large contribution to the reward like the size of the slate...
- Optimizing the policy is difficult...

Slate Recommendation (Framework)

From Single Items to Slates

- Given a policy π , the decision rule in single item recommendation is:

$$\operatorname{argmax}_a \pi(a | \mathbf{x})$$

- We want to reduce this decision rule to a MIPS: $\operatorname{argmax}_a f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_a$
- This is achievable using softmax: $\pi_{\Xi, \beta}(a | \mathbf{x}) = \frac{\exp(f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_a)}{\sum_{a'} \exp(f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_{a'})}$, in which case
$$\operatorname{argmax}_a \pi_{\Xi, \beta}(a | \mathbf{x}) = \operatorname{argmax}_a f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_a$$
- In slate recommendation, the decision rule becomes:

$$\operatorname{argmax}_{a_1, a_2, \dots, a_K} \pi(a_1, a_2, \dots, a_K | \mathbf{x})$$

- We want to reduce to a MIPS: $(\operatorname{argsort} f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta})_{1:K}$

Slate Recommendation (Framework)

From Single Items to Slates

- Given a policy π , the decision rule in single item recommendation is:

$$\operatorname{argmax}_a \pi(a | \mathbf{x})$$

- We want to reduce this decision rule to a MIPS: $\operatorname{argmax}_a f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_a$
- This is achievable using softmax: $\pi_{\Xi, \beta}(a | \mathbf{x}) = \frac{\exp(f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_a)}{\sum_{a'} \exp(f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_{a'})}$, in which case
$$\operatorname{argmax}_a \pi_{\Xi, \beta}(a | \mathbf{x}) = \operatorname{argmax}_a f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta}_a$$
- In slate recommendation, the decision rule becomes:

$$\operatorname{argmax}_{a_1, a_2, \dots, a_K} \pi(a_1, a_2, \dots, a_K | \mathbf{x})$$

- We want to reduce to a MIPS: $(\operatorname{argsort} f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta})_{1:K}$

The equivalence between argmax and argsort is not trivial.
The softmax parametrization does not solve it!

Slate Recommendation (Framework)

From Single Items to Slates

- In slate recommendation, the decision rule becomes:

$$\operatorname{argmax}_{a_1, a_2, \dots, a_K} \pi(a_1, a_2, \dots, a_K \mid \mathbf{x})$$

- We want to reduce to a MIPS: $\operatorname{argsort}(f_{\Xi}(\mathbf{x})^T \boldsymbol{\beta})_{1:K}$

→ Need of a clever parametrization of the models and policies...

The equivalence between
argmax and argsort is not trivial.
The softmax parametrization
does not solve it!



Slate

Recommendation

IPS Method

Slate Recommendation (IPS Method)

From Single Items to Slates

- Recall that the performance of recommendation policies π is often quantified by the value function:

$$V(\pi) = \mathbb{E}_{x \sim \nu} \mathbb{E}_{a_1, \dots, a_K \sim \pi(\cdot | x)} [R(x, a_1, \dots, a_K)] = \mathbb{E}_{x \sim \nu} \left[\sum_{a_1} \dots \sum_{a_K} \pi(a_1, \dots, a_K | x) R(x, a_1, \dots, a_K) \right]$$

Slate Recommendation (IPS Method)

From Single Items to Slates

- Recall that the performance of recommendation policies π is often quantified by the value function:

$$V(\pi) = \mathbb{E}_{x \sim \nu} \mathbb{E}_{a_1, \dots, a_K \sim \pi(\cdot | x)} [R(x, a_1, \dots, a_K)] = \mathbb{E}_{x \sim \nu} \left[\sum_{a_1} \dots \sum_{a_K} \pi(a_1, \dots, a_K | x) R(x, a_1, \dots, a_K) \right]$$

- The corresponding IPS estimator is

$$\hat{V}_n^{IPS}(\pi) = \frac{1}{n} \sum_{i=1}^n R_i \frac{\pi(a_{i,1}, \dots, a_{i,K} | x_i)}{\pi_0(a_{i,1}, \dots, a_{i,K} | x_i)}$$

Slate Recommendation (IPS Method)

From Single Items to Slates

- The IPS estimator of the value function is

$$\hat{V}_n^{IPS}(\pi) = \frac{1}{n} \sum_{i=1}^n R_i \frac{\pi(a_{i,1}, \dots, a_{i,K} | \mathbf{x}_i)}{\pi_0(a_{i,1}, \dots, a_{i,K} | \mathbf{x}_i)}$$

- It is unbiased.
- However, its variance scales linearly with $\frac{\pi(a_{i,1}, \dots, a_{i,K} | \mathbf{x}_i)}{\pi_0(a_{i,1}, \dots, a_{i,K} | \mathbf{x}_i)}$ which can be huge.
- Example:** The logging policy π_0 is uniform. In this case, the variance grows linearly with $P^K \pi(a | \mathbf{x})$, where P is the catalog size.

Slate Recommendation (IPS Method)

From Single Items to Slates

Solution: Independent IPS (IIPS)

- (A1) The expected reward of a slate is the sum of expected ranks of its items:

$$R(\mathbf{x}, a_1, \dots, a_K) = \sum_{k \in [K]} r_k(\mathbf{x}, a_1, \dots, a_K)$$

- (A2) The expected rank only depends on the item and its position:

$$r_k(\mathbf{x}, a_1, \dots, a_K) = r_k(\mathbf{x}, a_k)$$

- The Independent IPS (IIPS) estimator of the value function follows as

$$\hat{V}_n^{IIPS}(\pi) = \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K r_{i,k} \frac{\pi(a_{i,k}, k \mid \mathbf{x}_i)}{\pi_0(a_{i,k}, k \mid \mathbf{x}_i)}$$

Slate Recommendation (IPS Method)

From Single Items to Slates

Solution: Independent IPS (IIPS)

- (A1) The expected reward of a slate is the sum of expected ranks of its items:

$$R(\mathbf{x}, a_1, \dots, a_K) = \sum_{k \in [K]} r_k(\mathbf{x}, a_1, \dots, a_K)$$

- (A2) The expected rank only depends on the item and its position:

$$r_k(\mathbf{x}, a_1, \dots, a_K) = r_k(\mathbf{x}, a_k)$$

- The Independent IPS (IIPS) estimator of the value function follows as

$$\hat{V}_n^{IIPS}(\pi) = \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K r_{i,k} \frac{\pi(a_{i,k}, k \mid \mathbf{x}_i)}{\pi_0(a_{i,k}, k \mid \mathbf{x}_i)}$$

Marginal probability
that policy π places
item $a_{i,k}$ in position k

Slate Recommendation (IPS Method)

From Single Items to Slates

Solution: Independent IPS (IIPS)

- (A1) The expected reward of a slate is the sum of expected ranks of its items:

$$R(\mathbf{x}, a_1, \dots, a_K) = \sum_{k \in [K]} r_k(\mathbf{x}, a_1, \dots, a_K)$$

- (A2) The expected rank only depends on the item and its position:

$$r_k(\mathbf{x}, a_1, \dots, a_K) = r_k(\mathbf{x}, a_k)$$

- The Independent IPS (IIPS) estimator of the value function follows as

$$\hat{V}_n^{IIPS}(\pi) = \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K r_{i,k} \frac{\pi(a_{i,k}, k \mid \mathbf{x}_i)}{\pi_0(a_{i,k}, k \mid \mathbf{x}_i)}$$

Hard to compute

Slate Recommendation (IPS Method)

From Single Items to Slates

- The Independent IPS (IIPS) estimator of the value function is

$$\hat{V}_n^{IIPS}(\pi) = \frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K r_{i,k} \frac{\pi(a_{i,k}, k | \mathbf{x}_i)}{\pi_0(a_{i,k}, k | \mathbf{x}_i)}$$

- The IIPS estimator is unbiased under (A1) and (A2).
- The variance is reduced compared to IPS as it scales linearly with $\frac{\pi(a_{i,k}, k | \mathbf{x}_i)}{\pi_0(a_{i,k}, k | \mathbf{x}_i)}$ instead of $\frac{\pi(a_{i,1}, \dots, a_{i,K} | \mathbf{x}_i)}{\pi_0(a_{i,1}, \dots, a_{i,K} | \mathbf{x}_i)}$. Example: In the case of uniform logging policies, the variance is reduced from $P^K \pi(a | \mathbf{x})$ to $P\pi(a | \mathbf{x})$.

Slate Recommendation (IPS Method)

From Single Items to Slates

- IIIPS is unbiased and reduces the variance of IPS. However, IIIPS makes the strong independence assumption.
- Also, IIIPS requires computing the marginal distribution $\pi(a_{i,k}, k | \mathbf{x}_i)$, which is often intractable since

$$\pi(a', k | \mathbf{x}_i) = \underbrace{\sum_{a_1} \dots \sum_{a_K}}_{\text{This is totally infeasible}} \pi(a_1, \dots, a_K | \mathbf{x}) \mathbb{I}_{a_k=a'}$$

This is totally infeasible

Slate Recommendation (IPS Method)

From Single Items to Slates

- IIIPS is unbiased and reduces the variance of IPS. However, IIIPS makes the strong independence assumption.
- Also, IIIPS requires computing the marginal distribution $\pi(a_{i,k}, k \mid \mathbf{x}_i)$, which is often intractable since

$$\pi(a', k \mid \mathbf{x}_i) = \sum_{a_1} \dots \sum_{a_K} \pi(a_1, \dots, a_K \mid \mathbf{x}) \mathbb{I}_{a_k=a'}$$

Common Solution: Factorized policies

$$\pi(a_1, \dots, a_K \mid \mathbf{x}) = \prod_{k=1}^K \pi(a_k, k \mid \mathbf{x})$$

Slate Recommendation (IPS Method)

From Single Items to Slates

- IIIPS is unbiased and reduces the variance of IPS. However, IIIPS makes the strong independence assumption.
- Also, IIIPS requires computing the marginal distribution $\pi(a_{i,k}, k | \mathbf{x}_i)$, which is often intractable since

$$\pi(a', k | \mathbf{x}_i) = \sum_{a_1} \dots \sum_{a_K} \pi(a_1, \dots, a_K | \mathbf{x}) \mathbb{I}_{a_k=a'}$$

Common Solution: Factorized policies

$$\pi(a_1, \dots, a_K | \mathbf{x}) = \prod_{k=1}^K \pi(a_k, k | \mathbf{x})$$

This is the marginal probability that policy π places item a_k in position k

Slate Recommendation (IPS Method)

From Single Items to Slates

- IIIPS is unbiased and reduces the variance of IPS. However, IIIPS makes the strong independence assumption.
- Also, IIIPS requires computing the marginal distribution $\pi(a_{i,k}, k | \mathbf{x}_i)$, which is often intractable since $\pi(a', k | \mathbf{x}_i) = \sum_{a_1} \dots \sum_{a_K} \pi(a_1, \dots, a_K | \mathbf{x}) \mathbb{I}_{a_k=a'}$

Common Solution: Factorized policies

$$\pi_0(a_1, \dots, a_K | \mathbf{x}) = \prod_{k=1}^K \pi_0(a_k, k | \mathbf{x})$$

$$\pi(a_1, \dots, a_K | \mathbf{x}) = \prod_{k=1}^K \pi(a_k, k | \mathbf{x})$$

Problems:

- The logging policy π_0 is not necessarily factorized.
- Best slate by π contains the best item repeated K times.



Slate

Recommendation

Direct Method

Slate Recommendation (Direct Method)

Rank vs. Reward

- Ranking Model (e.g., BPR):
 - Trained on clicked slates only.
 - Not reward optimizing, it does not see it.
- Reward Model:
 - Trained on all slates.
 - Returns calibrated CTRs.



Slate Recommendation (Direct Method)

Rank vs. Reward

- Ranking Model (e.g., BPR):
 - Trained on clicked slates only.
 - Not reward optimizing, it does not see it.
- Reward Model:
 - Trained on all slates.
 - Returns calibrated CTRs.



WE CAN COMBINE BOTH → Probabilistic Reward and Rank (PRR [\(Aouali et al. 2022\)](#))

Slate Recommendation (Direct Method)

PRR Model ([Aouali et al. 2022](#))

The Probabilistic Rank and Reward Model (PRR) has the following properties:

- We assume at most one element in the slate is clicked.
 - This differs with IIPS, Position Based Model, etc.
- PRR includes user engagement features and recommendation features.
- PRR incorporates:
 - **was the slate clicked?** bandit approach
 - **which item was clicked?** IIPS, position-based models
- Decision making in PRR boils down to a MIPS task.

Details in: [Aouali et al. 2022](#)

Slate Recommendation (Direct Method)

PRR Model (Aouali et al. 2022)

$$\bar{R}, r_1, \dots, r_K | \mathbf{y}, \mathbf{z}, a_1, \dots, a_K \sim \text{categorical} \left(\frac{\theta_0}{Z}, \frac{\theta_1}{Z}, \dots, \frac{\theta_K}{Z} \right), \quad Z = \sum \theta_k$$

$\theta_0 = \exp(\mathbf{y}^T \boldsymbol{\phi})$

\mathbf{y} - User engagement features (slate level)

$\boldsymbol{\phi}$ - Nuisance parameters

Slate Recommendation (Direct Method)

PRR Model ([Aouali et al. 2022](#))

$$\bar{R}, r_1, \dots, r_K | \mathbf{y}, \mathbf{z}, a_1, \dots, a_K \sim \text{categorical} \left(\underbrace{\frac{\theta_0}{Z}, \frac{\theta_1}{Z}, \dots, \frac{\theta_K}{Z}}_{\sum \theta_k} \right), \quad Z = \sum \theta_k$$
$$\theta_k = \exp(g_\Gamma(\mathbf{z})^T \Psi_{a_k}) \exp(\gamma_k) + \exp(\alpha_k)$$

\mathbf{z} - context

$g_\Gamma(\mathbf{z})$ - contextual embedding

Ψ_a - Item embedding for item a

γ, α - multiplicative and additional position bias

Slate Recommendation (Direct Method)

PRR Learning ([Aouali et al. 2022](#))

$$\bar{R}, r_1, \dots, r_K | \mathbf{y}, \mathbf{z}, a_1, \dots, a_K \sim \text{categorical} \left(\frac{\theta_0}{Z}, \frac{\theta_1}{Z}, \dots, \frac{\theta_K}{Z} \right), \quad Z = \sum \theta_k$$
$$\theta_K = \exp(g_\Gamma(\mathbf{z})^T \Psi_{a_k}) \exp(\gamma_k) + \exp(\alpha_k)$$

\mathbf{z} - context

$g_\Gamma(\mathbf{z})$ - contextual embedding

Ψ_a - Item embedding for item a

γ, α - multiplicative and additional position bias

→ We use the maximum likelihood principle to estimate the parameters $\Psi, \Gamma, \gamma, \alpha, \phi$.

Slate Recommendation (Direct Method)

PRR Decision Rule ([Aouali et al. 2022](#))

$$\bar{R}, r_1, \dots, r_K | \mathbf{y}, \mathbf{z}, a_1, \dots, a_K \sim \text{categorical} \left(\frac{\theta_0}{Z}, \frac{\theta_1}{Z}, \dots, \frac{\theta_K}{Z} \right), \quad Z = \sum \theta_k$$
$$\rightarrow P(R = 1 | y, z, a_1, \dots, a_K) = 1 - \frac{\theta_0}{Z}$$

The decision rule follows as:

$$\operatorname{argmax}_{a_1, \dots, a_K} P(R = 1 | y, z, a_1, \dots, a_K) = \operatorname{argmax}_{a_1, \dots, a_K} 1 - \frac{\theta_0}{Z} = \operatorname{argmin}_{a_1, \dots, a_K} \frac{\theta_0}{Z}$$

Slate Recommendation (Direct Method)

PRR Decision Rule ([Aouali et al. 2022](#))

$$\bar{R}, r_1, \dots, r_K | \mathbf{y}, \mathbf{z}, a_1, \dots, a_K \sim \text{categorical} \left(\frac{\theta_0}{Z}, \frac{\theta_1}{Z}, \dots, \frac{\theta_K}{Z} \right), \quad Z = \sum \theta_k$$
$$\rightarrow P(R = 1 | y, z, a_1, \dots, a_K) = 1 - \frac{\theta_0}{Z}$$

The decision rule follows as:

$$\begin{aligned} \operatorname{argmax}_{a_1, \dots, a_K} P(R = 1 | y, z, a_1, \dots, a_K) &= \operatorname{argmax}_{a_1, \dots, a_K} 1 - \frac{\theta_0}{Z} = \operatorname{argmin}_{a_1, \dots, a_K} \frac{\theta_0}{Z} \\ &= \operatorname{argmin}_{a_1, \dots, a_K} \frac{\theta_0}{\theta_0 + \sum \theta_k} = \operatorname{argmax}_{a_1, \dots, a_K} \sum_{k=1}^K \theta_k \end{aligned}$$

Slate Recommendation (Direct Method)

PRR Decision Rule ([Aouali et al. 2022](#))

$$\bar{R}, r_1, \dots, r_K | \mathbf{y}, \mathbf{z}, a_1, \dots, a_K \sim \text{categorical} \left(\frac{\theta_0}{Z}, \frac{\theta_1}{Z}, \dots, \frac{\theta_K}{Z} \right), \quad Z = \sum \theta_k$$
$$\rightarrow P(R = 1 | y, z, a_1, \dots, a_K) = 1 - \frac{\theta_0}{Z}$$

The decision rule follows as:

$$\begin{aligned} \operatorname{argmax}_{a_1, \dots, a_K} P(R = 1 | y, z, a_1, \dots, a_K) &= \operatorname{argmax}_{a_1, \dots, a_K} 1 - \frac{\theta_0}{Z} = \operatorname{argmin}_{a_1, \dots, a_K} \frac{\theta_0}{Z} \\ &= \operatorname{argmin}_{a_1, \dots, a_K} \frac{\theta_0}{\theta_0 + \sum \theta_k} = \operatorname{argmax}_{a_1, \dots, a_K} \sum_{k=1}^K \theta_k \\ &= \operatorname{argmax}_{a_1, \dots, a_K} \sum_{k=1}^K \exp(g_\Gamma(\mathbf{z})^T \boldsymbol{\Psi}_{a_k}) \exp(\gamma_k) + \exp(\alpha_k) \end{aligned}$$

Slate Recommendation (Direct Method)

PRR Decision Rule ([Aouali et al. 2022](#))

$$\begin{aligned}\bar{R}, r_1, \dots, r_K | \mathbf{y}, \mathbf{z}, a_1, \dots, a_K &\sim \text{categorical} \left(\frac{\theta_0}{Z}, \frac{\theta_1}{Z}, \dots, \frac{\theta_K}{Z} \right), \quad Z = \sum \theta_k \\ \rightarrow P(R = 1 | y, z, a_1, \dots, a_K) &= 1 - \frac{\theta_0}{Z}\end{aligned}$$

The decision rule follows as:

$$\begin{aligned}\operatorname{argmax}_{a_1, \dots, a_K} P(R = 1 | y, z, a_1, \dots, a_K) &= \operatorname{argmax}_{a_1, \dots, a_K} 1 - \frac{\theta_0}{Z} = \operatorname{argmin}_{a_1, \dots, a_K} \frac{\theta_0}{Z} \\ &= \operatorname{argmin}_{a_1, \dots, a_K} \frac{\theta_0}{\theta_0 + \sum \theta_k} = \operatorname{argmax}_{a_1, \dots, a_K} \sum_{k=1}^K \theta_k \\ &= \operatorname{argmax}_{a_1, \dots, a_K} \sum_{k=1}^K \exp(g_\Gamma(\mathbf{z})^T \boldsymbol{\Psi}_{a_k}) \exp(\gamma_k) + \exp(\alpha_k) \\ &= \operatorname{argmax}_{a_1, \dots, a_K} \sum_{k=1}^K \exp(g_\Gamma(\mathbf{z})^T \boldsymbol{\Psi}_{a_k}) \exp(\gamma_k)\end{aligned}$$

Slate Recommendation (Direct Method)

PRR Decision Rule ([Aouali et al. 2022](#))

The decision rule follows as:

$$\operatorname{argmax}_{a_1, \dots, a_K} P(R = 1 \mid y, z, a_1, \dots, a_K) = \operatorname{argmax}_{a_1, \dots, a_K} \sum_{k=1}^K \exp(g_\Gamma(z)^T \Psi_{a_k}) \exp(\gamma_k)$$

The argmax operation on the right-hand-side can be solved using MIPS as follows:

$$a'_1, \dots, a'_K = \operatorname{argsort}(g_\Gamma(z)^T \Psi)_{1:K}$$

$$i_1, \dots, i_K = \operatorname{argsort}(\gamma)$$

$$a_1, \dots, a_K = a'_{i_1}, \dots, a'_{i_K}$$

Slate Recommendation (Direct Method)

PRR Variants (Aouali et al. 2022)

PRR model

$$\bar{R}, r_1, \dots, r_K | \mathbf{y}, \mathbf{z}, a_1, \dots, a_K \sim \text{categorical} \left(\frac{\theta_0}{Z}, \frac{\theta_1}{Z}, \dots, \frac{\theta_K}{Z} \right), \quad Z = \sum \theta_k$$

PRR Reward Only

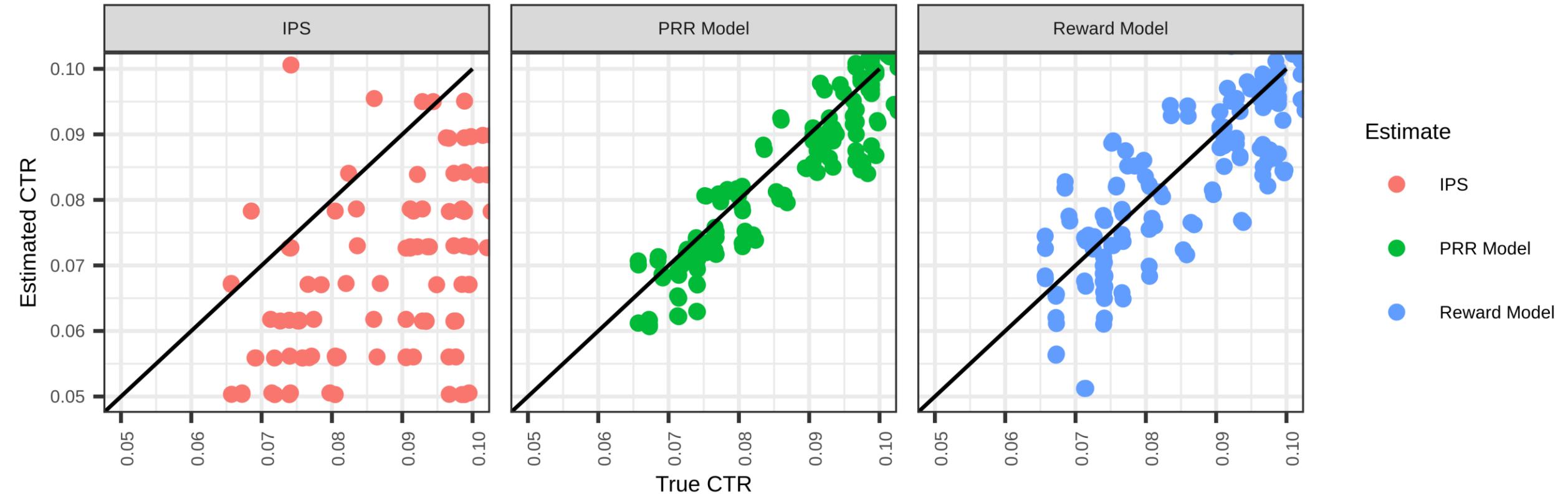
$$\bar{R}, R | \mathbf{y}, \mathbf{z}, a_1, \dots, a_K \sim \text{categorical} \left(\frac{\theta_0}{Z}, \frac{\sum \theta_k}{Z} \right), \quad Z = \sum_{k=0}^K \theta_k$$

PRR Rank Only

$$r_1, \dots, r_K | R = 1, \mathbf{y}, \mathbf{z}, a_1, \dots, a_K \sim \text{categorical} \left(\frac{\theta_1}{Z}, \dots, \frac{\theta_K}{Z} \right), \quad Z = \sum_{k=1}^K \theta_k$$

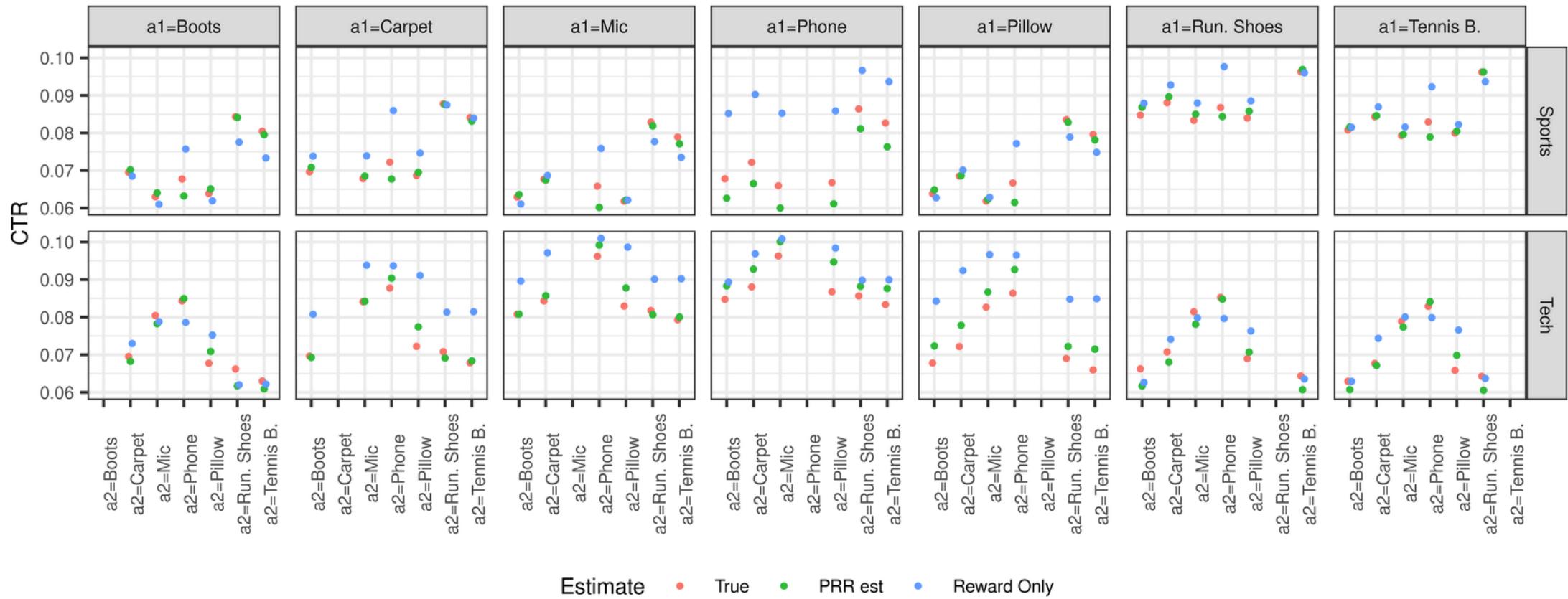
$$\theta_0 = \exp(\mathbf{y}^T \boldsymbol{\phi}) \text{ and } \theta_k = \exp(g_\Gamma(\mathbf{z})^T \boldsymbol{\Psi}_{ak}) \exp(\gamma_k) + \exp(\alpha_k)$$

Slate Recommendation (Direct Method)



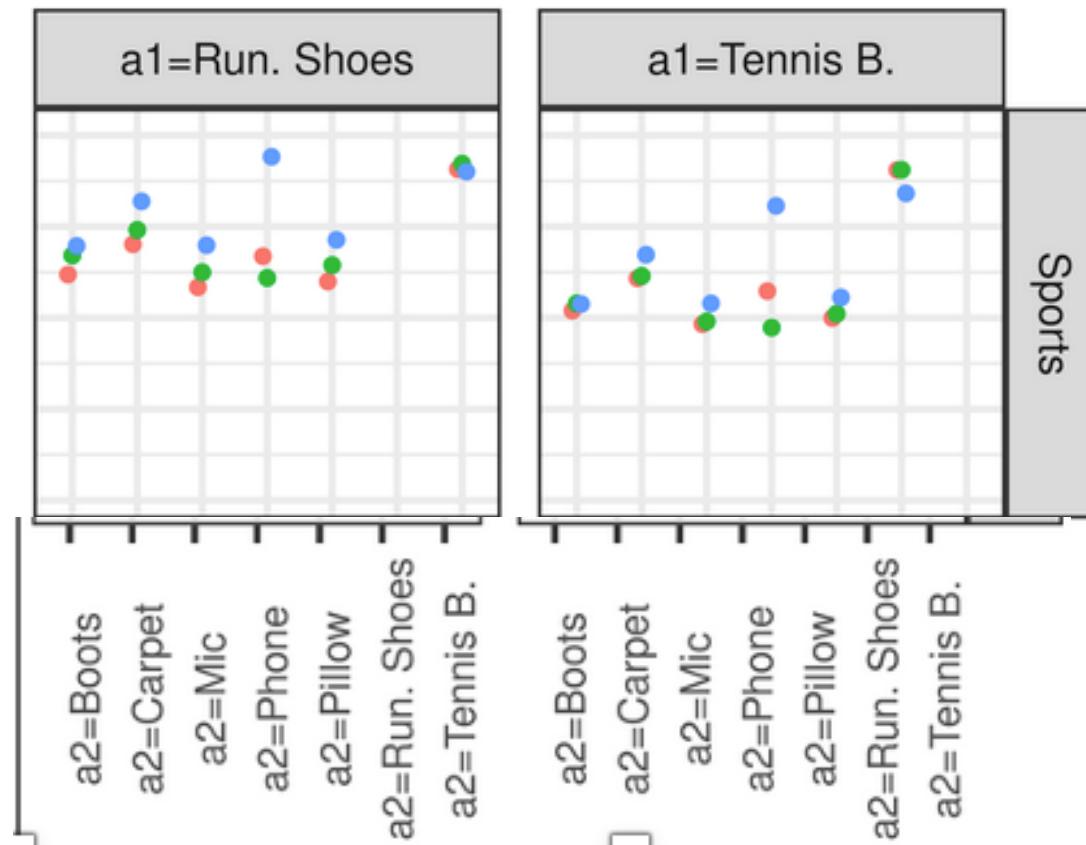
Slate Recommendation (Direct Method)

PRR Toy Example: CTR Estimates

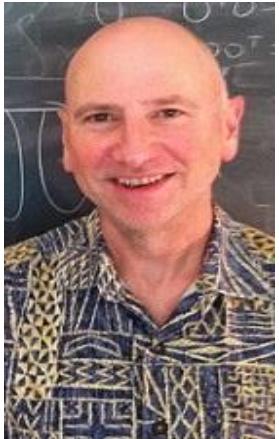


Slate Recommendation (Direct Method)

PRR Toy Example: CTR Estimates



Slate Recommendation (DM vs. IPS)

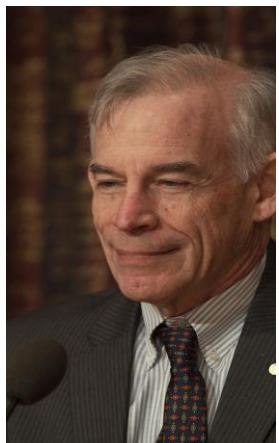


- Christopher Sims and Larry Wasserman had an epic debate on the IPS (or Horvitz-Thompson) estimation: Whether it causes a problem for likelihood/Bayesian inference.
- Bayesians/Likelihoodists should not use the propensity score...
- PRR + BLOB ignore the propensity score and outperform other methods.

Debate and discussion:

[Bayesian Causal Inference for Real World Interactive Systems \(bcirwis2021.github.io\)](https://bcirwis2021.github.io)

[Christopher Sims - Large Parameter Spaces and Weighted Data: A Bayesian Perspective - YouTube](https://www.youtube.com/watch?v=KJLjyfzXWUw)





Slate

Recommendation

Decision Rule

Optimization

Slate Recommendation (Decision Rule Optimization)

Policy Learning for Slates

$$\hat{V}_{\Xi,\beta}(x) = \sum_{a_1} \dots \sum_{a_K} \hat{R}(x, a_1, \dots, a_K) \pi_{\Xi,\beta}(a_1, \dots, a_K | x)$$

$$\nabla_{\Xi,\beta} \hat{V}_{\Xi,\beta}(x) = \nabla_{\Xi,\beta} \sum_{a_1} \dots \sum_{a_K} \hat{R}(x, a_1, \dots, a_K) \pi_{\Xi,\beta}(a_1, \dots, a_K | x)$$

This is now totally infeasible

Slate Recommendation (Decision Rule Optimization)

Reinforce Algorithm for Slates

$$\hat{V}_{\Xi,\beta}(x) = \sum_{a_1} \dots \sum_{a_K} \hat{R}(x, a_1, \dots, a_K) \pi_{\Xi,\beta}(a_1, \dots, a_K | x)$$

Solution:

$$a'_1, \dots, a'_K \sim \pi_{\Xi,\beta}(\cdot | x)$$

$$\nabla_{\Xi,\beta} \hat{V}_{\Xi,\beta}(x) = \hat{R}(x, a'_1, \dots, a'_K) \nabla_{\Xi,\beta} \log \pi_{\Xi,\beta}(a'_1, \dots, a'_K | x)$$

If we can sample and evaluate $\pi_{\Xi,\beta}(\cdot | x)$ then we can proceed...

Slate Recommendation (Decision Rule Optimization)

Reinforce Algorithm for Slates

$$\hat{V}_{\Xi,\beta}(x) = \sum_{a_1} \dots \sum_{a_K} \hat{R}(x, a_1, \dots, a_K) \pi_{\Xi,\beta}(a_1, \dots, a_K | x)$$

Solution:

$$a'_1, \dots, a'_K \sim \pi_{\Xi,\beta}(\cdot | x)$$

$$\nabla_{\Xi,\beta} \hat{V}_{\Xi,\beta}(x) = \hat{R}(x, a'_1, \dots, a'_K) \nabla_{\Xi,\beta} \log \pi_{\Xi,\beta}(a'_1, \dots, a'_K | x)$$

If we can sample and evaluate $\pi_{\Xi,\beta}(\cdot | x)$ then we can proceed...



Factorized policies

Slate Recommendation (Decision Rule Optimization)

Back to Factorized Policies



$$\pi(a_1, \dots, a_K | x) = \prod_k \pi(a_k, k | x)$$

What if the best two recommendations are rice and beer. Suppose that

$$\begin{aligned}\pi(a = \text{rice} | x) &= 0.51 \\ \pi(a = \text{beer} | x) &= 0.49\end{aligned}$$

Then

$$\begin{aligned}\pi(a_1 = \text{rice}, a_2 = \text{rice} | x) &= 0.2601 \\ \pi(a_1 = \text{rice}, a_2 = \text{beer} | x) &= 0.2499 \\ \pi(a_1 = \text{beer}, a_2 = \text{rice} | x) &= 0.2499 \\ \pi(a_1 = \text{beer}, a_2 = \text{beer} | x) &= 0.2401\end{aligned}$$

Slate Recommendation (Decision Rule Optimization)

Back to Factorized Policies

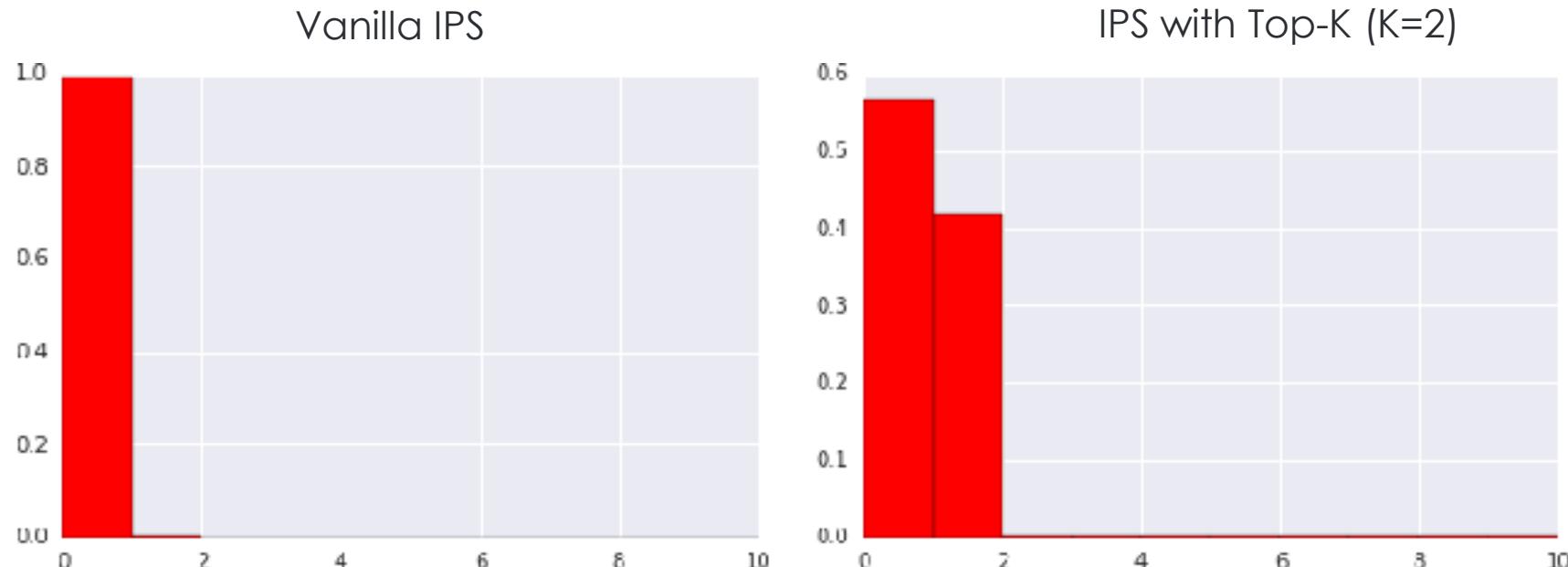


$$\pi(a_1, \dots, a_K | x) = \prod_k \pi(a_k, k | x)$$

Solution: Chen, et al. 2019 proposes a solution to this which keeps the probability spread out over the Top-K items.

Slate Recommendation (Decision Rule Optimization)

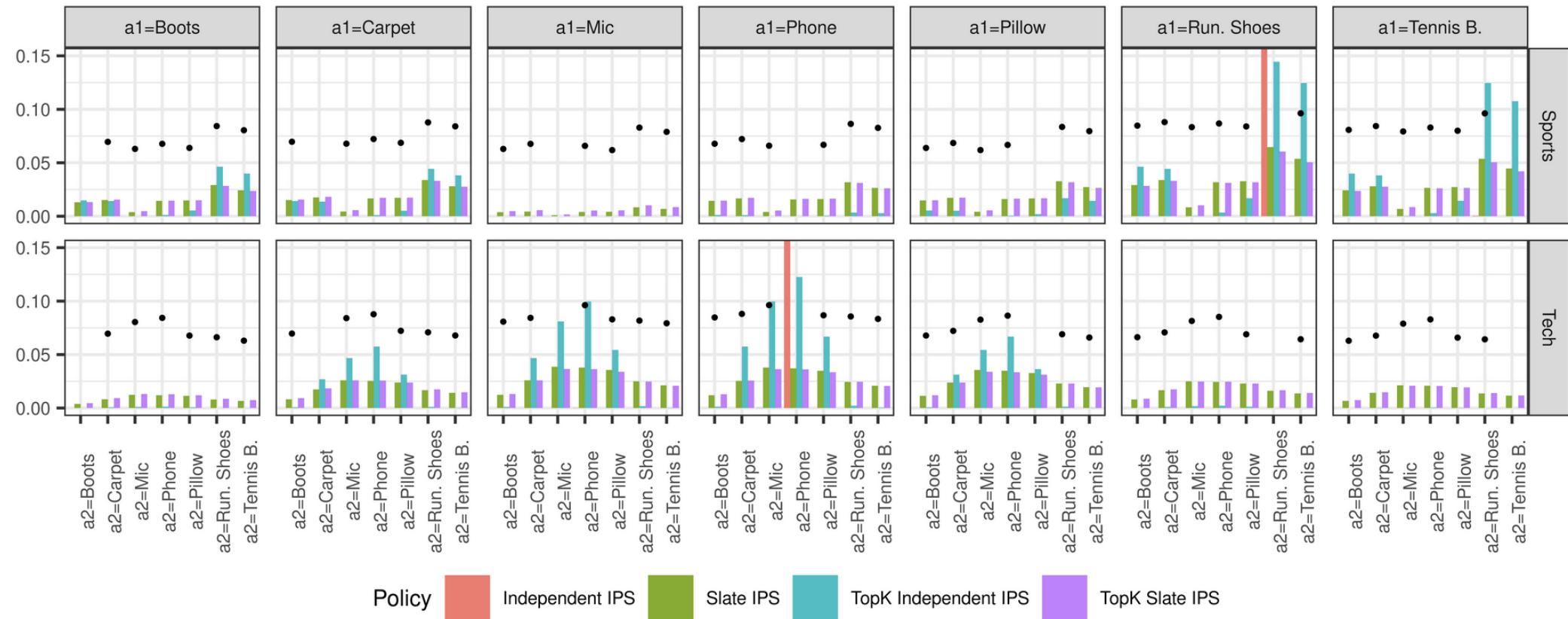
Factorized Policies with Top-K



The Top-K correction modifies the loss to stop $\pi(a_k, k|x)$ converging to a degenerate distribution.

Slate Recommendation (Decision Rule Optimization)

Toy Example: IPS vs. IIIPS



Slate Recommendation (Practice)

Python Tutorial 3

https://colab.research.google.com/github/otmhi/Reward-Optimizing-Reco/blob/main/ECML_Reward_Optimizing_Slate_Recommendation_with_DL_and_MIP_S_Part_3.ipynb

Takeaways

Takeaways

- DL + MIPS is powerful - scales to large catalogues and can fit near arbitrary models.
- A/B tests produce a population split and do an excellent job of estimating reward. In contrast off-policy estimation either with DM or IPS are difficult to apply in practice:
 - Contextual bandit assumption is typically made (reward is on the recommendations not the timeline).
 - IPS: Very high variance even in idealized circumstances.
 - Practitioners fall back on proxy methods which mean that the deep learning objective and A/B test results are not aligned.

Takeaways

This tutorial surveys approaches to improving this alignment in a practical MIPS setting:

- Decision theoretic framework introduced, separating reward estimation from optimizing the decision rule.
- Methods for improving the reward estimation on full slates include IIPS and the direct method PRR.
- Some parametric models result in a MIPS compatible decision rule automatically.
- When this is not possible or desirable, optimizing the decision rule can be done using policy learning or the reinforce algorithm.
- A softmax policy is suitable for single item slates, a factorized policy can work for multiple item slates but can draw duplicates.

Takeaways

This tutorial surveys approaches to improving this alignment in a practical MIPS setting:

- Successfully achieving positive A/B tests requires:
 - Efficient reward estimation of recommendations.
 - Efficient optimization of MIPS decision rule.

References

Aouali, Imad, et al. "Probabilistic Rank and Reward: A Scalable Model for Slate Recommendation." *arXiv preprint arXiv:2208.06263* (2022).

Chen, Minmin, et al. "Top-k off-policy correction for a REINFORCE recommender system." Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining. 2019.

Beygelzimer, Alina, et al. "Contextual bandit algorithms with supervised learning guarantees." *Proceedings of the Fourteenth International Conference on Artificial Intelligence and Statistics*. JMLR Workshop and Conference Proceedings, 2011.

Bottou, Léon, et al. "Counterfactual Reasoning and Learning Systems: The Example of Computational Advertising." *Journal of Machine Learning Research* 14.11 (2013).

Ding, Qin, Hsiang-Fu Yu, and Cho-Jui Hsieh. "A fast sampling algorithm for maximum inner product search." *The 22nd International Conference on Artificial Intelligence and Statistics*. PMLR, 2019.

References

Dudík, Miroslav, et al. "Doubly robust policy evaluation and optimization." *Statistical Science* 29.4 (2014): 485-511.

Joachims, Thorsten, Adith Swaminathan, and Maarten De Rijke. "Deep learning with logged bandit feedback." *International Conference on Learning Representations*. 2018.

Li, Shuai, et al. "Offline evaluation of ranking policies with click models." *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2018.

Malkov, Yu A., and Dmitry A. Yashunin. "Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs." *IEEE transactions on pattern analysis and machine intelligence* 42.4 (2018): 824-836.

McInerney, James, et al. "Counterfactual evaluation of slate recommendations with sequential reward interactions." *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2020.

References

Robins, James M., and Ya'acov Ritov. "Toward a curse of dimensionality appropriate (CODA) asymptotic theory for semi-parametric models." *Statistics in medicine* 16.3 (1997): 285-319.

Saito, Yuta, and Thorsten Joachims. "Counterfactual learning and evaluation for recommender systems: Foundations, implementations, and recent advances." *Fifteenth ACM Conference on Recommender Systems*. 2021.

Sakhi, Otmane, et al. "BLOB: A probabilistic model for recommendation that combines organic and bandit signals." *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 2020.

Sakhi, Otmane, David Rohde, and Alexandre Gilotte. "Fast Offline Policy Optimization for Large Scale Recommendation." arXiv preprint arXiv:2208.05327 (2022).

Swaminathan, Adith, and Thorsten Joachims. "Batch learning from logged bandit feedback through counterfactual risk minimization." *The Journal of Machine Learning Research* 16.1 (2015): 1731-1755.